

Machine Learning

Volker Roth

Department of Mathematics & Computer Science
University of Basel

Section 4

Regression

Regression basics

- In regression we assume that a response variable $y \in \mathbb{R}$ is a noisy function of the input variable $\mathbf{x} \in \mathbb{R}^d$.

$$y = f(\mathbf{x}) + \eta.$$

- We often assume that f is linear, $f(\mathbf{x}) = \mathbf{w}^t \mathbf{x}$, and that η has a zero-mean Gaussian distribution with constant variance, $\eta \sim N(0, \sigma^2)$.
- This can equivalently be written as

$$p(y|\mathbf{x}) = N(\mu(\mathbf{x}), \sigma^2), \text{ with } \mu(\mathbf{x}) = \mathbf{w}^t \mathbf{x}.$$

- In one dimension: $\mu(\mathbf{x}) = w_0 + w_1 x$ and $\mathbf{x} = (1, x)$.
 w_0 is the **intercept** or bias term and w_1 is the **slope**.
- If $w_1 > 0$, we expect the output to increase as the input increases.

Least Squares and Maximum Likelihood

- Fit n data points (\mathbf{x}_i, y_i) to a model that has $d + 1$ parameters w_j , $j = 0, \dots, d$.
- Notation: $\mathbf{x} \leftarrow (1, \mathbf{x}) \rightsquigarrow w_0$ is the intercept.
- Frequentist view: \mathbf{w} is an unknown parameter vector, not a RV.
- We assume that the n observations are **iid**.
- **Linear model:** $y_i = \mathbf{w}^t \mathbf{x}_i + \eta_i$, $\eta_i \sim N(0, \sigma^2)$.
Observed y_i generated from a normal distribution centered at $\mathbf{w}^t \mathbf{x}_i$.
- Model predicts linear relationship between **conditional expectation** of **observations** y_i and **inputs** \mathbf{x}_i :

$$E[y_i | \mathbf{x}_i] = w_0 + w_1 x_{i1} + \dots + w_d x_{id} = \mathbf{w}^t \mathbf{x}_i = f(\mathbf{x}_i; \mathbf{w}).$$

Note: the expectation operator is linear and $E[\eta_i] = 0$.

Regression function = conditional expectation.

LS and Maximum Likelihood

- **Likelihood function:** conditional probability of all observed y_i given their explanation, treated as a function of the model parameters \mathbf{w} :

$$L(\mathbf{w}) \propto \prod_i \exp \left[-\frac{1}{2\sigma^2} (y_i - \mathbf{w}^t \mathbf{x}_i)^2 \right]$$

- **Maximizing L = finding model that best explains observations:**

$$\begin{aligned} \hat{\mathbf{w}} &= \arg \max_{\mathbf{w}} L(\mathbf{w}) = \arg \min_{\mathbf{w}} [-L(\mathbf{w})] = \arg \min_{\mathbf{w}} [-\log(L(\mathbf{w}))] \\ &= \arg \min_{\mathbf{w}} \sum_i (y_i - \mathbf{w}^t \mathbf{x}_i)^2 \end{aligned}$$

Least-squares fit = ML solution under Gaussian error model.

- $\hat{\mathbf{w}}_{MLE}$ minimizes the **residual sum of squares**

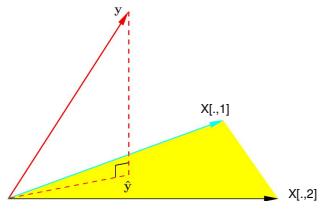
$$RSS(\mathbf{w}) = \sum_{i=1}^n r_i^2 = \sum_{i=1}^n [y_i - f(\mathbf{x}_i; \mathbf{w})]^2 = \|\mathbf{y} - X\mathbf{w}\|^2.$$

Least squares regression: Geometry

$$\begin{aligned}\frac{\partial RSS(\mathbf{w})}{\partial \mathbf{w}} &= \frac{\partial}{\partial \mathbf{w}} [\mathbf{y}^t \mathbf{y} - 2\mathbf{y}^t X \mathbf{w} + \mathbf{w}^t X^t X \mathbf{w}] \\ &= -2X^t \mathbf{y} + 2X^t X \mathbf{w} \stackrel{!}{=} \mathbf{0} \\ \Rightarrow \quad \hat{\mathbf{w}} &= (X^t X)^{-1} X^t \mathbf{y} \\ \Rightarrow \quad X^t (\mathbf{y} - X \hat{\mathbf{w}}) &= X^t \hat{\mathbf{r}} = \mathbf{0}.\end{aligned}$$

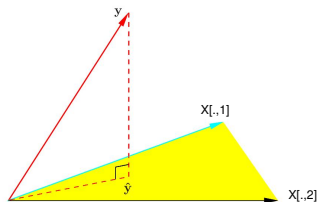
It follows that $\sum_{i=1}^n X_{ij} r_i = 0$, $\forall j = 0, 1, \dots, d$.

Residual is orthogonal to 1 ($j = 0$) and to every input dimension $X_{\bullet j}$.



Adapted from Fig. 3.2 in (Hastie, Tibshirani, Friedman)

Least squares regression: Geometry



Adapted from Fig. 3.2 in (Hastie, Tibshirani, Friedman)

- The fitted values at the training inputs are

$$(\hat{f}(\mathbf{x}_1), \dots, \hat{f}(\mathbf{x}_n))^t = \hat{\mathbf{y}} = \mathbf{X}\hat{\mathbf{w}} = \mathbf{X}(\mathbf{X}^t\mathbf{X})^{-1}\mathbf{X}^t\mathbf{y}.$$

- $H = \mathbf{X}(\mathbf{X}^t\mathbf{X})^{-1}\mathbf{X}^t$ is called “hat” matrix (puts hat on \mathbf{y})
- Column vectors of \mathbf{X} span the **column space** of $\mathbf{X} \subset \mathbb{R}^n$.
- Minimizing $RSS(\mathbf{w}) \rightsquigarrow$ choose $\hat{\mathbf{w}}$ such that \mathbf{r} is orthogonal.
- Fitted values $\hat{\mathbf{y}}$ are **orthogonal projection** of \mathbf{y} on column space.

Least squares regression: Algebra

- H is **orthogonal projection** on **column space** of X :

$$HX = X(X^tX)^{-1}X^tX = X.$$

- **Fundamental theorem of linear algebra:** the **nullspace** of X^t is the orthogonal complement of the column space of X .

- $M = I_n - H$ is **orthogonal projection** on **nullspace** of X^t :

$$MX = (I_n - H)X = X - X = 0.$$

- H and M are symmetric ($H^t = H$) and idempotent ($MM = M$)

The Algebra of Least Squares

- H creates fitted values: $\hat{\mathbf{y}} = H\mathbf{y} \rightsquigarrow \hat{\mathbf{y}} \in \text{Col}(X)$
- M creates residuals: $\mathbf{r} = M\mathbf{y} \rightsquigarrow \hat{\mathbf{r}} \in \text{Null}(X^t) \Leftrightarrow X^t\mathbf{r} = \mathbf{0}$

Frequentist confidence limits

- **Recall:** $y_i = f(\mathbf{x}_i; \mathbf{w}) + \eta_i$, with independent Gaussian noise.
- In matrix-vector form: $\mathbf{y} = X\mathbf{w} + \boldsymbol{\eta}$, with $\boldsymbol{\eta} \sim N(\mathbf{0}, \sigma^2 I_n)$.

$$\begin{aligned}\hat{\mathbf{w}} &= (X^t X)^{-1} X^t \mathbf{y} \\ &= (X^t X)^{-1} X^t X \mathbf{w} + (X^t X)^{-1} X^t \boldsymbol{\eta} \\ &= \mathbf{w} + (X^t X)^{-1} X^t \boldsymbol{\eta}\end{aligned}$$

$$\Rightarrow \quad \hat{\mathbf{w}} - \mathbf{w} = (X^t X)^{-1} X^t \boldsymbol{\eta} =: A\boldsymbol{\eta}$$

- **Linear** functions of normals are normal:

$$\boldsymbol{\eta} \sim N(\mathbf{0}, \sigma^2 I_n) \Rightarrow A\boldsymbol{\eta} \sim N(\mathbf{0}, \sigma^2 A A^t).$$

$$\text{Here: } A = (X^t X)^{-1} X^t \Rightarrow A A^t = (X^t X)^{-1}$$

- Conditioned on X and σ^2 :

$$\hat{\mathbf{w}} - \mathbf{w} | X, \sigma^2 \sim N\left(\mathbf{0}, \sigma^2 (X^t X)^{-1}\right).$$

Frequentist confidence limits

- Distribution completely specified \leadsto **confidence limits:**

$$\hat{w}_k - w_k \sim N(0, \sigma^2 S^{kk}),$$

where S^{kk} denotes the k th diagonal element of $(X^t X)^{-1}$.

- Thus, both z'_k and $z_k = -z'_k$ are standard normal:

$$z_k := (w_k - \hat{w}_k) / \sqrt{\sigma^2 S^{kk}} \sim N(0, 1)$$

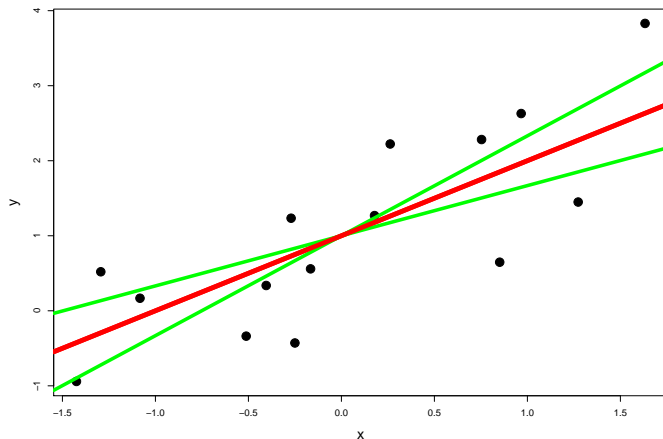
- CDF:

$$P(z_k < k_c) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{k_c} e^{-t^2/2} dt =: \Phi(k_c) = 1 - c$$

- Upper limit for w_k :

$$\begin{aligned} P(z_k < k_c) &= P(\sqrt{\sigma^2 S^{kk}} z_k < \sqrt{\sigma^2 S^{kk}} k_c) \\ &= P(w_k - (w_k - \hat{w}_k) > w_k - \sqrt{\sigma^2 S^{kk}} k_c) \\ &= P(\hat{w}_k > w_k - \sqrt{\sigma^2 S^{kk}} k_c) \\ &= P(w_k < \hat{w}_k + \sqrt{\sigma^2 S^{kk}} k_c) = 1 - c. \end{aligned}$$

Frequentist confidence limits



Least-squares fit (red) and two lines with slopes according to upper (lower) 95% confidence limit (green).

Standard parametric rate

- Assume we have estimated the parameters based on n samples:

$$\begin{aligned}(\hat{\mathbf{w}}_n - \mathbf{w}) &\sim N(\mathbf{0}, \sigma^2 (X^t X)^{-1}) \\&= N(\mathbf{0}, \sigma^2 (X^t X / n)^{-1} \cdot 1/n) \\ \sqrt{n}(\hat{\mathbf{w}}_n - \mathbf{w}) &\sim N(\mathbf{0}, \sigma^2 \underbrace{(X^t X / n)^{-1}}_{\rightarrow \Sigma})\end{aligned}$$

- Since for $n \rightarrow \infty$, $X^t X / n \rightarrow \Sigma = \text{const}$, this means that **$\hat{\mathbf{w}}_n$ converges to \mathbf{w} at a rate of $1/\sqrt{n}$.**
- This is a very general result that holds in an asymptotic sense even without assuming normality \rightsquigarrow **central limit theorem.**
- Due to its universality, it is called the **standard parametric rate.**

Basis functions

- Can be generalized to model non-linear relationships by replacing \mathbf{x} with some non-linear function of the inputs, $\phi(\mathbf{x})$:

$$p(y|\mathbf{x}) = N(\mathbf{w}^t \phi(\mathbf{x}), \sigma^2).$$

- Predictions can be based on a linear combination of a set of basis functions $\phi(\mathbf{x}) = \{g_0(\mathbf{x}), g_1(\mathbf{x}), \dots, g_m(\mathbf{x})\}$, with $g_i(\mathbf{x}) : \mathbb{R}^d \mapsto \mathbb{R}$. Can model the intercept by setting $g_0(\mathbf{x}) = 1$:

$$f(\mathbf{x}; \mathbf{w}) = w_0 + w_1 g_1(\mathbf{x}) + \dots + w_m g_m(\mathbf{x}).$$

\leadsto **additive models**

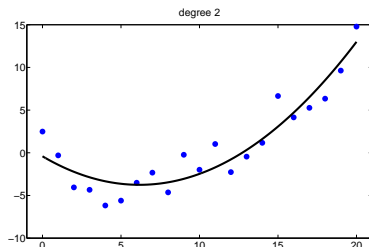
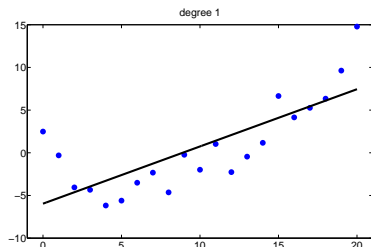


Fig 1.7 in K.Murphy

Additive models

- Examples:

If $x \in \mathbb{R}^d$ and $m = d + 1$, $g_0(\mathbf{x}) = 1$ and $g_i(\mathbf{x}) = x_i, i = 1, \dots, d$, then

$$f(\mathbf{x}; \mathbf{w}) = w_0 + w_1 x_1 + \dots + w_d x_d.$$

If $x \in \mathbb{R}$, $g_0(\mathbf{x}) = 1$ and $g_i(x) = x^i, i = 1, \dots, m$, then

$$f(x; \mathbf{w}) = w_0 + w_1 x^1 + \dots + w_m x^m.$$

- Basis functions can **capture various properties of the inputs**.

Example: **Document analysis**

\mathbf{x} = text document (collection of words)

$$g_i(\mathbf{x}) = \begin{cases} 1, & \text{if word } i \text{ appears in the document} \\ 0, & \text{otherwise} \end{cases}$$

$$f(\mathbf{x}; \mathbf{w}) = w_0 + \sum_{i \in \text{words}} w_i g_i(\mathbf{x}).$$

Additive models cont'd

- We can also make predictions by gauging the **similarity of examples to prototypes**.
- For example, our additive regression function could be

$$f(\mathbf{x}; \mathbf{w}) = w_0 + w_1 g_1(\mathbf{x}) + \cdots + w_m g_m(\mathbf{x}),$$

where the basis functions are **radial basis functions**

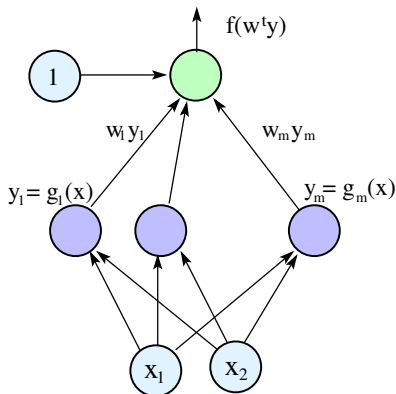
$$g_k(\mathbf{x}) = \exp\left(-\frac{1}{2\sigma^2} \|\mathbf{x} - \mathbf{x}_k\|^2\right)$$

measuring the similarity to the prototypes \mathbf{x}_k .

- The variance σ^2 controls how quickly the basis function vanishes as a function of the distance to the prototype.
- **Training examples themselves could serve as prototypes.**

Additive models cont'd

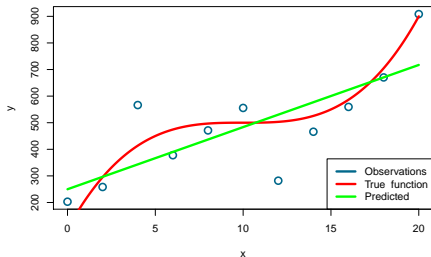
Can view additive models graphically in terms of **units** and **weights**.



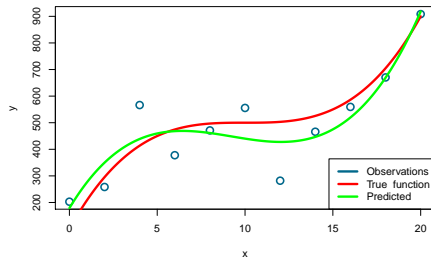
In **neural networks** the basis functions have adjustable parameters.

Example: Polynomial regression

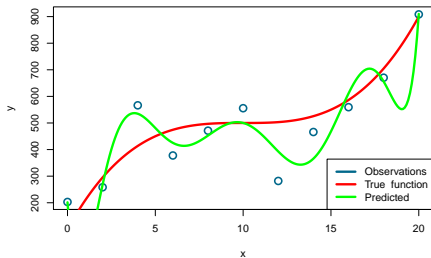
Polynomial basis functions. Degree = 1



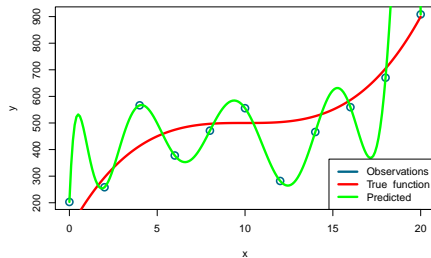
Polynomial basis functions. Degree = 3



Polynomial basis functions. Degree = 8



Polynomial basis functions. Degree = 10



Complexity and overfitting

With limited training examples our polynomial regression model may achieve zero training error but nevertheless has a large expected error.

$$\text{training} \quad \frac{1}{n} \sum_{i=1}^n (y_i - f(\mathbf{x}_i; \hat{\mathbf{w}}))^2 \approx 0$$

$$\text{expectation} \quad E_{(\mathbf{x}, y) \sim p} (y - f(\mathbf{x}; \hat{\mathbf{w}}))^2 \gg 0$$

We suffer from **over-fitting**

\rightsquigarrow should reconsider our model \rightsquigarrow **model selection.**

We will discuss model selection from a **Bayesian perspective** first.

A frequentist approach will follow later in the chapter on **statistical learning theory.**

Subsection 1

Bayesian Regression

Bayesian interpretation: priors

- Suppose our generative model takes an input $\mathbf{x} \in \mathbb{R}^d$ and maps it to a real valued output y according to

$$p(y|\mathbf{x}, \mathbf{w}, \sigma^2) = N(y|\mathbf{w}^t \mathbf{x}, \sigma^2)$$

- We will keep σ^2 fixed and only try to estimate \mathbf{w} .
- Given data $\mathcal{D} = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$, the **likelihood function** is

$$L(\mathbf{w}; \mathcal{D}) = \prod_{i=1}^n N(y_i|\mathbf{w}^t \mathbf{x}_i, \sigma^2) = \prod_{i=1}^n \frac{1}{Z} \exp\left(-\frac{1}{2\sigma^2}(y_i - \mathbf{w}^t \mathbf{x}_i)^2\right).$$

- In classical regression we used the maximizing parameters $\hat{\mathbf{w}}$.
- **In Bayesian analysis we keep all regression functions**, just weighted by their ability to explain the data.
- Our knowledge about \mathbf{w} after seeing the data is defined by the **posterior distribution** $p(\mathbf{w}|\mathcal{D})$.

Bayesian regression: Prior and posterior

- We specify our **prior belief** about the parameter values as $p(\mathbf{w})$.
For instance, we could prefer small parameter values:

$$p(\mathbf{w}) = N(\mathbf{w} | 0, \tau^2 I)$$

The smaller τ^2 is, the smaller values of \mathbf{w} we prefer
prior to seeing the data.

- **Posterior** proportional to prior $p(\mathbf{w})$ times likelihood:

$$p(\mathbf{w} | \mathcal{D}) \propto L(\mathbf{w}; \mathcal{D}) p(\mathbf{w})$$

- Here: posterior is **Gaussian** $p(\mathbf{w} | \mathcal{D}, \sigma^2) = N(\mathbf{w} | \mathbf{w}_N, V_N)$ with mean \mathbf{w}_N and covariance V_N given by

$$\mathbf{w}_N = (X^t X + \lambda I)^{-1} X^t \mathbf{y}, \quad V_N = \sigma^2 (X^t X + \lambda I)^{-1},$$

with $\lambda = \frac{\sigma^2}{\tau^2}$.

Bayesian regression: Posterior computation

Given variables $\mathbf{x} \in \mathbb{R}^{d_x}$ and $\mathbf{y} \in \mathbb{R}^{d_y}$, assume **linear Gaussian system**:

$$p(\mathbf{x}) = N(\mathbf{x} | \boldsymbol{\mu}_x, \Sigma_x) \quad (\leadsto \text{prior})$$

$$p(\mathbf{y} | \mathbf{x}) = N(\mathbf{y} | A\mathbf{x} + \mathbf{b}, \Sigma_y) \quad (\leadsto \text{likelihood})$$

- **The posterior is also Gaussian:**

$$p(\mathbf{x} | \mathbf{y}) = N(\mathbf{x} | \boldsymbol{\mu}_{x|y}, \Sigma_{x|y})$$

$$\Sigma_{x|y}^{-1} = \Sigma_x^{-1} + A^t \Sigma_y^{-1} A$$

$$\boldsymbol{\mu}_{x|y} = \Sigma_{x|y} \left(A^t \Sigma_y^{-1} (\mathbf{y} - \mathbf{b}) + \Sigma_x^{-1} \boldsymbol{\mu}_x \right).$$

Gaussian likelihood and Gaussian prior form a conjugate pair.

- The normalization constant (denominator in Bayes formula) is

$$p(\mathbf{y}) = N(\mathbf{y} | A\boldsymbol{\mu}_x + \mathbf{b}, \Sigma_y + A\Sigma_x A^t).$$

Bayesian regression: Posterior predictive

- Prediction of y for new \mathbf{x} : use posterior as weights for predictions based on individual \mathbf{w} 's \rightsquigarrow **Posterior predictive:**

$$\begin{aligned} p(y|\mathbf{x}, \mathcal{D}, \sigma^2) &= \int p(y|\mathbf{x}, \mathbf{w}, \sigma^2) p(\mathbf{w}|\mathcal{D}) d\mathbf{w} \\ &= \int N(y|\mathbf{x}^t \mathbf{w}, \sigma^2) N(\mathbf{w}|\mathbf{w}_N, V_N) \\ &= N(y|\mathbf{w}_N^t \mathbf{x}, \sigma_N^2(\mathbf{x})), \text{ with} \\ \sigma_N^2(\mathbf{x}) &= \sigma^2 + \mathbf{x}^t V_N \mathbf{x}. \end{aligned}$$

- The variance in this prediction, $\sigma_N^2(\mathbf{x})$, depends on two terms:
 - ▶ the variance of the observation noise, σ^2
 - ▶ the variance in the parameters, V_N
 - \rightsquigarrow depends on how close \mathbf{x} is to training data \mathcal{D}
 - \rightsquigarrow error bars get larger as we move away from training points.

Bayesian regression: Posterior predictive

- By contrast, the **plugin approximation** uses only the ML-parameter estimate with the degenerate distribution $p(\mathbf{w}|\mathcal{D}, \sigma^2) = \delta_{\hat{\mathbf{w}}}(\mathbf{w})$:
$$p(\mathbf{y}|\mathbf{x}, \mathcal{D}, \sigma^2) \approx \int p(\mathbf{y}|\mathbf{x}, \mathbf{w}, \sigma^2) \delta_{\hat{\mathbf{w}}}(\mathbf{w}) d\mathbf{w} = p(\mathbf{y}|\mathbf{x}, \hat{\mathbf{w}}, \sigma^2) = N(\mathbf{y}|\mathbf{x}^t \hat{\mathbf{w}}, \sigma^2).$$

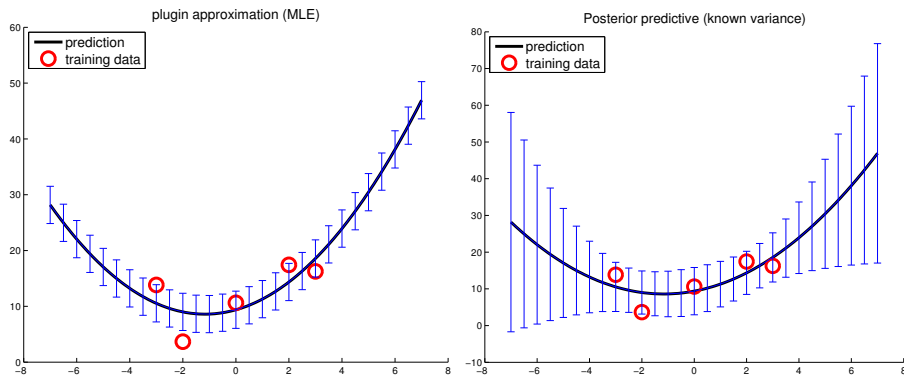


Fig. 7.12 in (K. Murphy). Example with quadratic basis functions: posterior predictive distribution (mean and $\pm 1\sigma$).

Sampling from posterior predictive

Left: plugin approximation: $f(y) = \phi(\mathbf{x})^t \hat{\mathbf{w}}$,
where $\phi(\mathbf{x})$ is the expanded input vector $(1, x, x^2)^t$.

Right: sampled functions $\phi(\mathbf{x})^t \mathbf{w}^{(s)}$, where $w^{(s)}$ are samples from the posterior

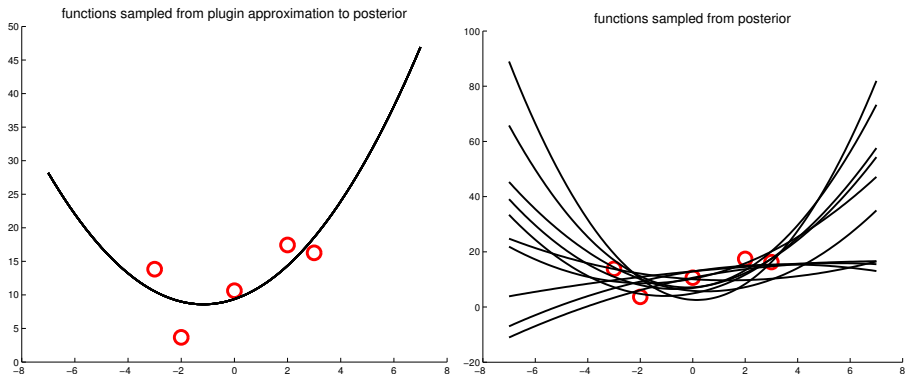


Fig. 7.12 in (K. Murphy)

MAP approximation and ridge regression

- Posterior proportional to prior $p(\mathbf{w}) = N(\mathbf{w}|0, \tau^2 I)$ times likelihood.
- The MAP estimate is

$$\begin{aligned}\mathbf{w}_{\text{MAP}} &= \arg \max \{ \log[L(\mathbf{w}; \mathcal{D})] + \log[p(\mathbf{w})] \} \\ &= \arg \min \{ -\log[L(\mathbf{w}; \mathcal{D})] - \log[p(\mathbf{w})] \} \\ &= \arg \min \left\{ \frac{1}{2\sigma^2} \sum_i (y_i - \mathbf{w}^t \mathbf{x}_i)^2 + \frac{1}{2\tau^2} \mathbf{w}^t \mathbf{w} \right\} \\ &= \arg \min \left\{ \sum_i (y_i - \mathbf{w}^t \mathbf{x}_i)^2 + \frac{\sigma^2}{\tau^2} \mathbf{w}^t \mathbf{w} \right\} \\ &= \arg \min \left\{ \sum_i (y_i - \mathbf{w}^t \mathbf{x}_i)^2 + \lambda \mathbf{w}^t \mathbf{w} \right\}\end{aligned}$$

- In classical statistics, this is called **ridge regression**:

$$\mathbf{w}_{\text{MAP}} = \mathbf{w}_{\text{ridge}} = (X^t X + \lambda I)^{-1} X^t \mathbf{y}.$$

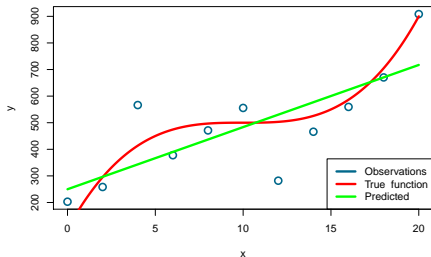
- In regularization theory, this is an example of **Tikhonov Regularization**.

Subsection 2

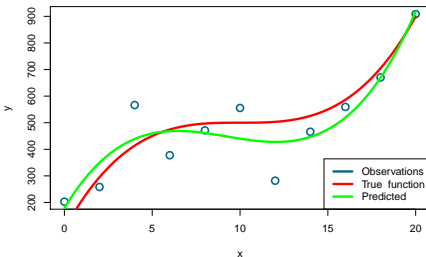
Bayesian model selection

Example: Polynomial regression

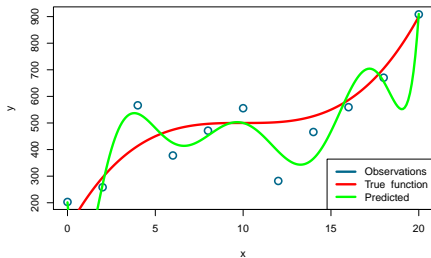
Polynomial basis functions. Degree = 1



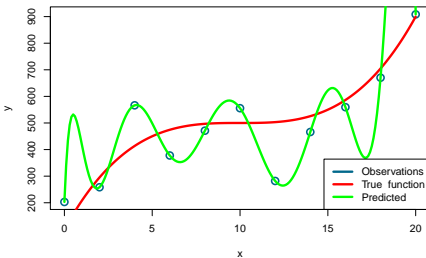
Polynomial basis functions. Degree = 3



Polynomial basis functions. Degree = 8



Polynomial basis functions. Degree = 10



Bayesian regression (again)

- Suppose our parametrized model \mathcal{F}_θ takes an input $\mathbf{x} \in \mathbb{R}^d$ and maps it to a real valued output y according to

$$p(y|\mathbf{x}, \boldsymbol{\theta}, \sigma^2) = N(y; \boldsymbol{\theta}^t \mathbf{x}, \sigma^2)$$

- We will keep σ^2 fixed and only try to estimate $\boldsymbol{\theta}$.
- Given data $\mathcal{D} = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$, define likelihood

$$L(\boldsymbol{\theta}; \mathcal{D}) = \prod_{i=1}^n N(y_i; \boldsymbol{\theta}^t \mathbf{x}_i, \sigma^2) = \prod_{i=1}^n \frac{1}{Z} \exp \left(-\frac{1}{2\sigma^2} (y_i - \boldsymbol{\theta}^t \mathbf{x}_i)^2 \right).$$

- In classical regression we used the maximizing parameters $\hat{\boldsymbol{\theta}}$.
- In Bayesian analysis we keep **all regression functions**, just weighted by their ability to explain the data.
- Knowledge about $\boldsymbol{\theta}$ after seeing the data defined by posterior $p(\boldsymbol{\theta}|\mathcal{D})$.

Bayesian regression (again)

- We specify our **prior belief** about the parameter values as $p(\theta)$.
For instance, we could prefer small parameter values:

$$p(\theta) = N(\theta; 0, \tau^2 I)$$

Small $\tau^2 \rightsquigarrow$ small θ preferred **prior to seeing data**.

- Posterior proportional to prior $p(\theta)$ times likelihood:

$$p(\theta|\mathcal{D}) \propto L(\theta; \mathcal{D})p(\theta)$$

- Normalization constant, a.k.a. **marginal likelihood**:

$$p(\mathbf{y}|\mathcal{F}, X) = \int \underbrace{L(\theta; \mathcal{D})}_{p(\mathbf{y}|\theta, X)} p(\theta|\mathcal{F}) d\theta,$$

depends on model + data but **not on specific parameter values**.

Example: Bayesian regression

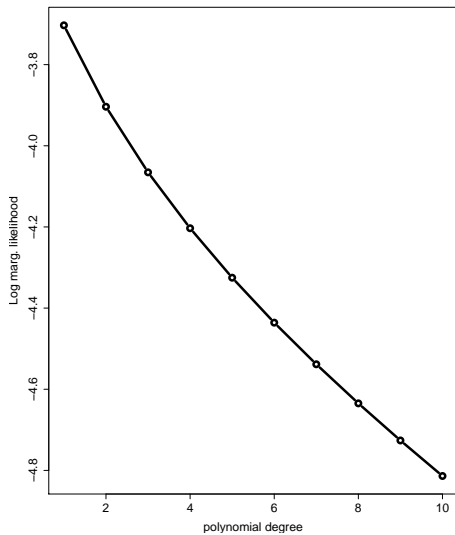
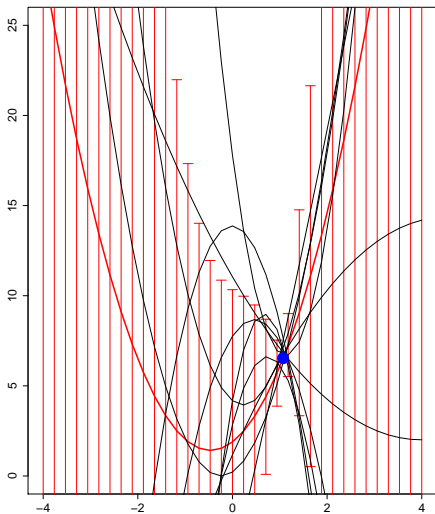
- Goal: choose among regression model families, specified by different feature mappings $\mathbf{x} \rightarrow \phi(\mathbf{x})$.
- Example: linear $\phi_1(\mathbf{x})$ and quadratic $\phi_2(\mathbf{x})$.

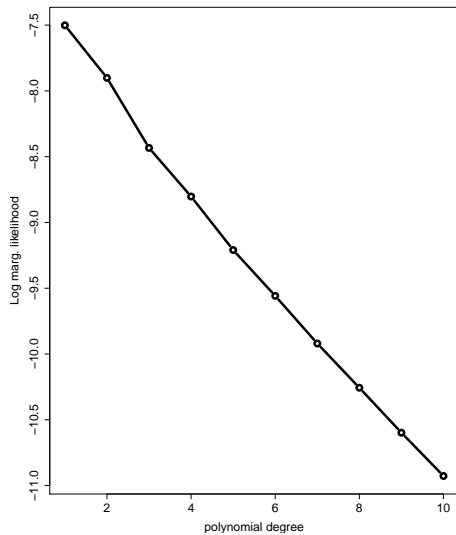
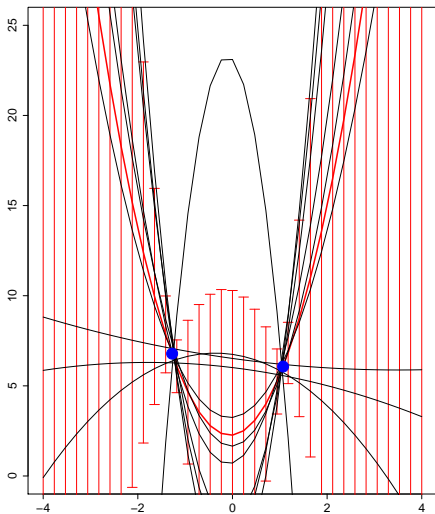
- The model families we compare are:

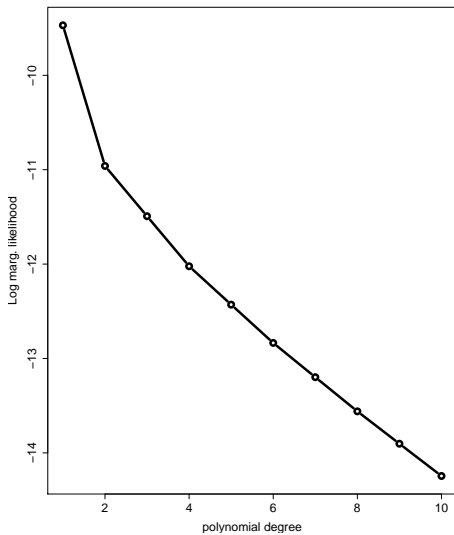
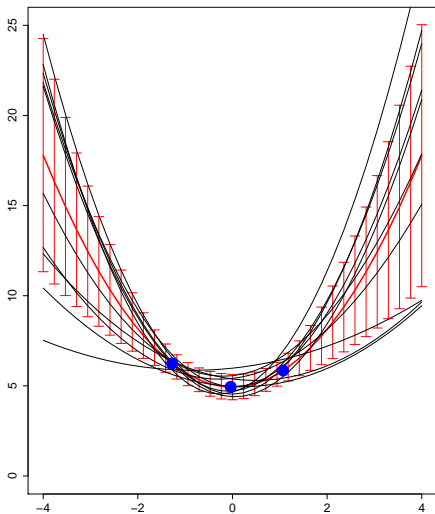
$$\mathcal{F}_1 : p(\mathbf{y}|\mathbf{x}, \theta_1, \sigma^2) = N(\mathbf{y}|\theta_1^t \phi_1(\mathbf{x}), \sigma^2)$$

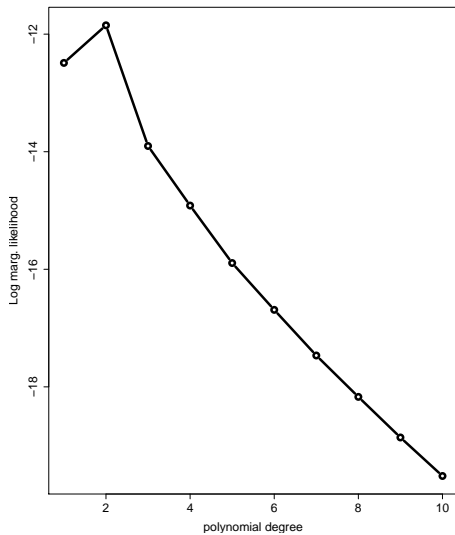
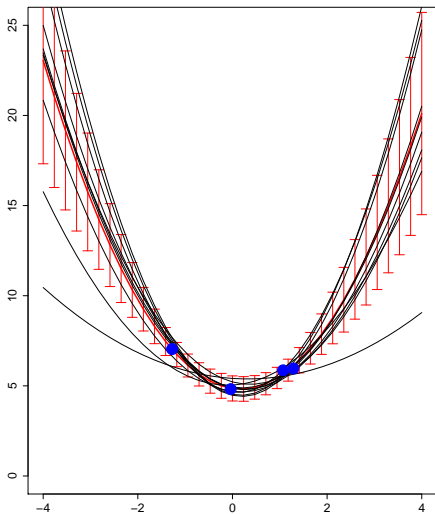
$$\mathcal{F}_2 : p(\mathbf{y}|\mathbf{x}, \theta_2, \sigma^2) = N(\mathbf{y}|\theta_2^t \phi_2(\mathbf{x}), \sigma^2).$$

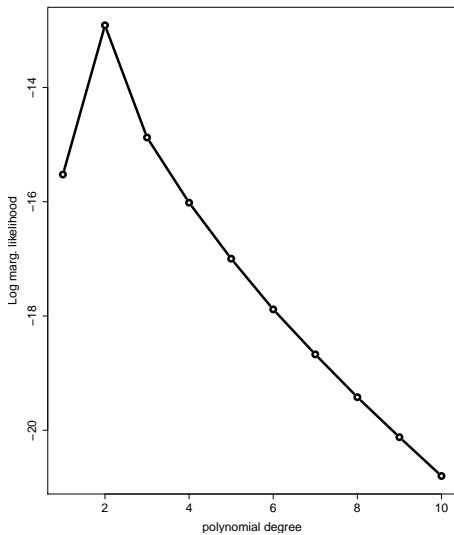
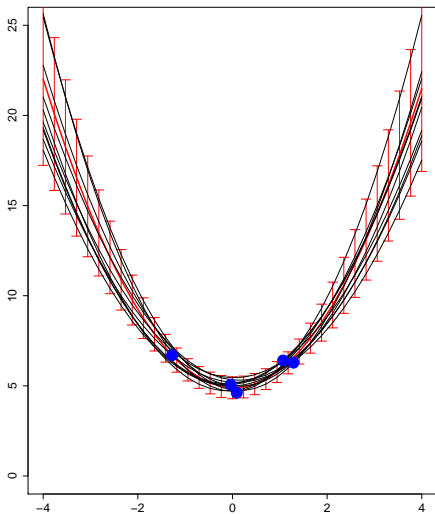
- Focusing on $p(\mathbf{y}|\mathcal{F}, X) = \int L(\theta; \mathcal{D})p(\theta)d\theta$, two possibilities:
 - ▶ **\mathcal{F} too flexible:** posterior $p(\theta|\mathcal{D})$ requires many training examples before it focuses on useful parameter values;
 - ▶ **\mathcal{F} too simple:** posterior concentrates quickly but the predictions remain poor.
- Pragmatic choice: Select the family whose **marginal likelihood** (a.k.a. **Bayesian score**) is larger.
- After seeing data \mathcal{D} we would select model \mathcal{F}_1 if $p(\mathbf{y}|\mathcal{F}_1, X) > p(\mathbf{y}|\mathcal{F}_2, X)$.

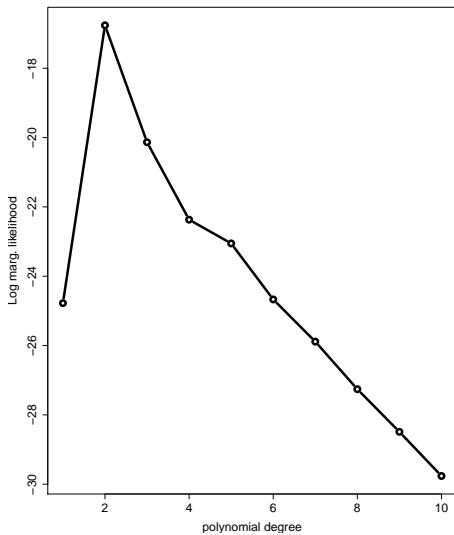
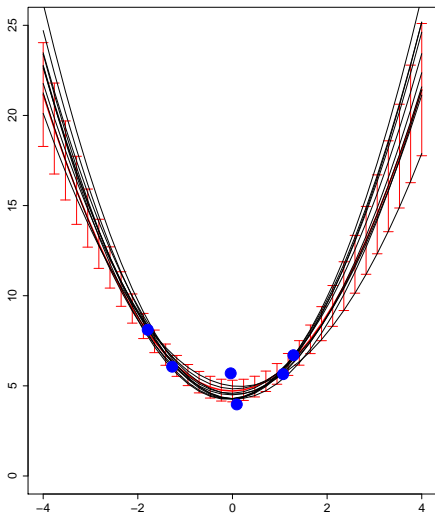


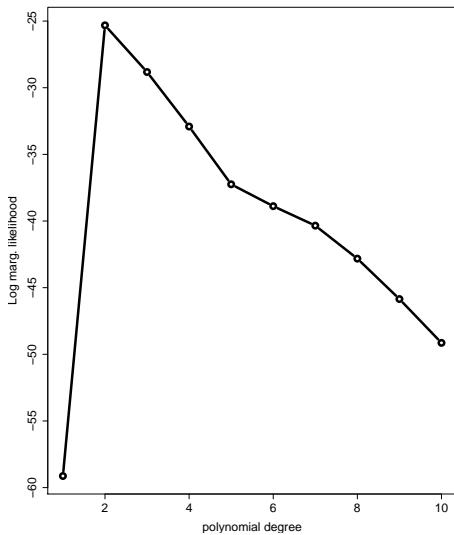
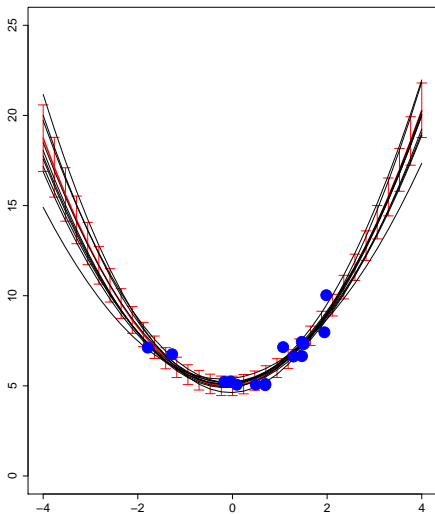












Approximating the marginal likelihood

- Problem: In most cases we cannot compute the marginal likelihood in closed form \rightsquigarrow approximations are needed.
- A specific approximation will lead to the **Bayesian Information Criterion (BIC)**.
- Key insight: when computing

$$p(\mathbf{y}|\mathcal{F}, X) = \int p(\mathbf{y}|\boldsymbol{\theta}, X)p(\boldsymbol{\theta}|\mathcal{F})d\boldsymbol{\theta},$$

the integrand is a product of two densities \rightsquigarrow integrand itself is an unnormalized density.

- Laplace's approximation uses a clever trick to approximate such integrals...

Approximation details: Laplace's Method

- Assume unnormalized density $p^*(\theta)$ has peak at $\hat{\theta}$. Goal: calculate normalizing constant

$$Z_p = \int p^*(\theta) d\theta$$

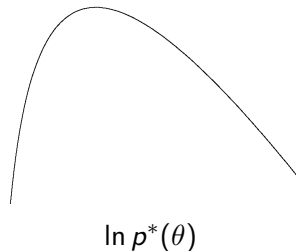
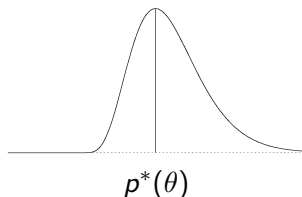
- Taylor-expand logarithm around $\hat{\theta}$:

$$\ln p^*(\theta) \approx \ln p^*(\hat{\theta}) - \frac{c}{2}(\theta - \hat{\theta})^2 + \dots,$$

where

$$c := -\frac{\partial^2}{\partial \theta^2} \ln p^*(\theta) \Big|_{\theta=\hat{\theta}}.$$

(note that first order term vanishes)



Laplace's Method (cont'd)

- Approximate $p^*(\theta)$ by unnormalized Gaussian

$$Q^*(\theta) := p^*(\hat{\theta}) \exp \left[-c/2 \cdot (\theta - \hat{\theta})^2 \right]$$

- A normalized Gaussian would be:

$$Q(\theta \mid \mu = \hat{\theta}, \sigma^2) = \frac{1}{Z_Q} \exp \left[-\frac{(\theta - \hat{\theta})^2}{2\sigma^2} \right],$$

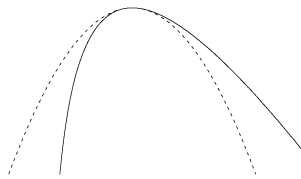
$$\text{with } Z_Q = \sqrt{2\pi\sigma^2} = \int \exp \left[-\frac{(\theta - \hat{\theta})^2}{2\sigma^2} \right] d\theta$$

- Approximate $Z_p = \int p^*(\theta) d\theta$ by

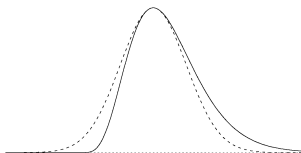
$$Z_p \approx \int Q^*(\theta) d\theta$$

$$= p^*(\hat{\theta}) \int \exp \left[-c/2 \cdot (\theta - \hat{\theta})^2 \right] d\theta$$

$$= p^*(\hat{\theta}) \sqrt{2\pi/c} \rightsquigarrow c \text{ is the inverse variance}$$



$\ln p^*(\theta)$ & $\ln Q^*(\theta)$



$p^*(\theta)$ & $Q^*(\theta)$

Laplace's Method (cont'd)

- Multivariate generalization in d dimensions:
second derivative \rightsquigarrow **Hessian matrix**

$$H_{ij} = \left. \frac{\partial^2 \ln p^*(\theta)}{\partial \theta_i \partial \theta_j} \right|_{\theta = \hat{\theta}}$$
$$Z_p \approx p^*(\hat{\theta}) \int \exp \left[-\frac{1}{2} (\theta - \hat{\theta})^t H (\theta - \hat{\theta}) \right] d\theta$$
$$= p^*(\hat{\theta}) \sqrt{\frac{(2\pi)^d}{|H|}} = p^*(\hat{\theta}) \left| \frac{H}{2\pi} \right|^{-\frac{1}{2}},$$

where the last equation follows from the properties of the determinant: $|aM| = a^d |M|$ for $M \in \mathbb{R}^{d \times d}$, $a \in \mathbb{R}$.

- Another interpretation: complicated distribution $p(\theta)$ is approximated by Gaussian centered at the mode $\hat{\theta}$:

$$p(\theta) \approx \mathcal{N}(\theta | \mu = \hat{\theta}, \Sigma = H^{-1}).$$

Example: Bayesian logistic regression

- Linear logistic regression: model parameters are simply the weights \mathbf{w} .
- Likelihood: $p(y|\mathbf{x}, \mathbf{w}) = \text{Ber}(y|\text{sigm}(\mathbf{w}^t \mathbf{x}))$
- Unfortunately, there is no convenient conjugate prior. Let's use a standard Gaussian prior: $p(\mathbf{w}) = N(\mathbf{w}|\mathbf{0}, V_0)$
- Laplace's approximation of posterior:

$$p(\mathbf{w}|\mathcal{D}) \approx N(\mathbf{w}|\mathbf{w}^*, H^{-1})$$

$$\mathbf{w}^* = \arg \max J[\mathbf{w}], \quad J[\mathbf{w}] = \log \underbrace{p(y|\mathbf{x}, \mathbf{w})}_{\text{likelihood}} + \log \underbrace{p(\mathbf{w})}_{\text{prior}}$$

$$H = \nabla^2 J(\mathbf{w}) \Big|_{\mathbf{w}^*}$$

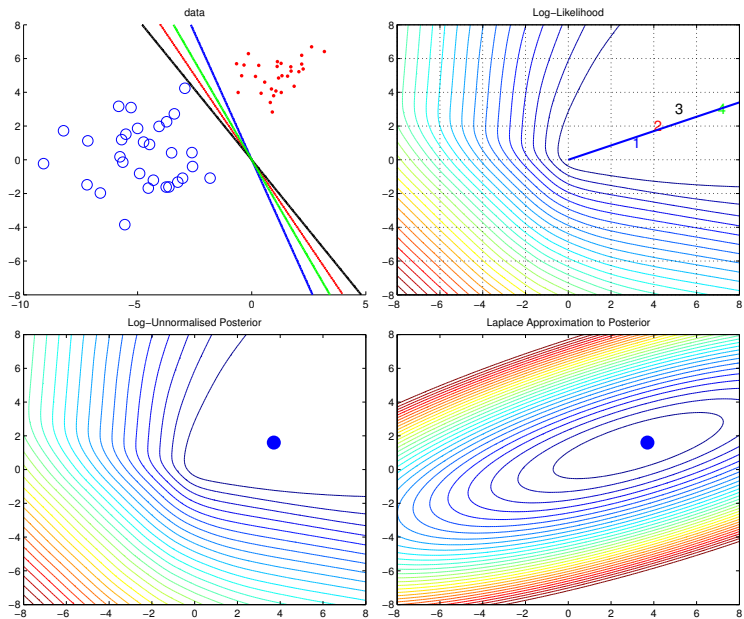


Fig 8.5 in K.Murphy

Bayesian LOGREG: Approximating the posterior predictive

- Posterior \rightsquigarrow can compute credible intervals etc.
- But in machine learning, interest usually focuses on **prediction**.
- The **posterior predictive distribution** has the form

$$p(y|\mathbf{x}, \mathcal{D}) = \int p(y|\mathbf{x}, \mathbf{w})p(\mathbf{w}|\mathcal{D}) d\mathbf{w}.$$

Here (and in most cases), this integral is intractable.

- The simplest approximation is the plug-in approximation

$$p(y = 1|\mathbf{x}, \mathcal{D}) \approx p(y = 1|\mathbf{x}, \mathbf{w}^*)$$

- But such a plug-in estimate underestimates the uncertainty.
- Better: **Monte Carlo approximation**

$$p(y|\mathbf{x}, \mathcal{D}) \approx \frac{1}{S} \sum_{s=1}^S \text{sigm}((\mathbf{w}^s)^t \mathbf{x}),$$

where $\mathbf{w}^s \sim p(\mathbf{w}|\mathcal{D})$ are samples from the Gaussian approximation to the posterior.

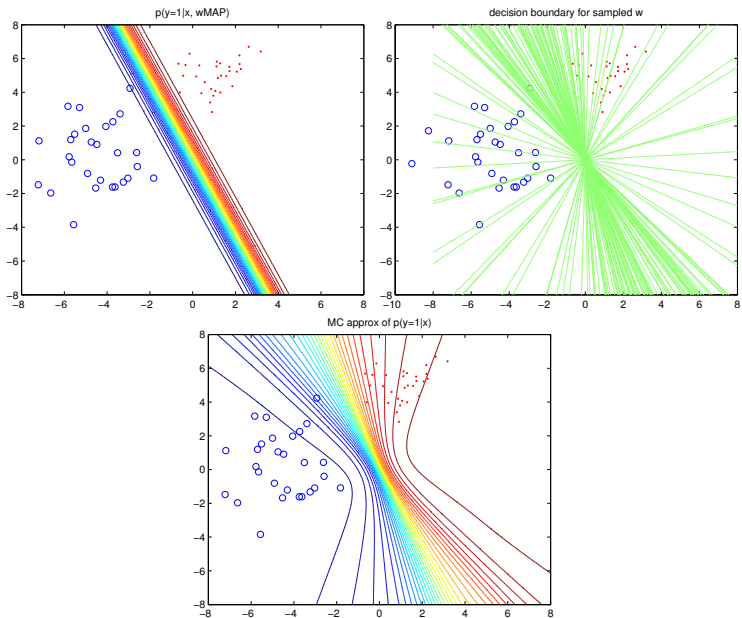


Fig 8.6 in K.Murphy

Approximating the marginal likelihood

$$\begin{aligned}p(\mathcal{D}|\mathcal{F}) &= \int p(\mathcal{D}|\boldsymbol{\theta}) \cdot p(\boldsymbol{\theta}|\mathcal{F}) d\boldsymbol{\theta} \\&\approx p(\mathcal{D}|\boldsymbol{\theta}^*) \cdot p(\boldsymbol{\theta}^*|\mathcal{F}) |H/(2\pi)|^{-\frac{1}{2}} \overset{\text{flat prior}}{\approx} p(\mathcal{D}|\hat{\boldsymbol{\theta}}) |H/(2\pi)|^{-\frac{1}{2}} \\ \log p(\mathcal{D}|\mathcal{F}) &\approx \log p(\mathcal{D}|\hat{\boldsymbol{\theta}}) - \frac{1}{2} \log |H| + C, \quad \text{with } \hat{\boldsymbol{\theta}} = \boldsymbol{\theta}_{MLE} \text{ in } \mathcal{F}.\end{aligned}$$

- Focus on last term:

$$H = \sum_{i=1}^n H_i, \quad \text{with } H_i = \nabla_{\boldsymbol{\theta}} \nabla_{\boldsymbol{\theta}} \log p(\mathcal{D}_i|\boldsymbol{\theta}).$$

Let's approximate each H_i with a **fixed** matrix H'

$$\log |H| = \log |nH'| = \log(n^d |H'|) = d \log n + \log(|H'|).$$

- For model selection, last term can be dropped, because it is independent of \mathcal{F} and n .

$$\log p(\mathcal{D}|\mathcal{F}) \approx \log p(\mathcal{D}|\hat{\boldsymbol{\theta}}) - \frac{d}{2} \log n + C = \text{BIC}(\mathcal{F}, n|\mathcal{D}) + C.$$

Intuitive interpretation of BIC

- The **Shannon information content** of a specific outcome a of a random experiment is

$$h(a) = -\log_2 P(a) = \log \frac{1}{P(a)}.$$

It measures the “surprise” (in bits):

Outcomes that are less probable have larger values of surprise.

- **Information theory:** Can find a code so that the **number of bits** used to encode each symbol $a \in \mathcal{A}$ is essentially $-\log_2 P(a)$.
- Here:

$$-\text{BIC}(\mathcal{F}, n|\mathcal{D}) = \overbrace{\sum_{i=1}^n \left(\underbrace{-\log_2 p(y_i | \mathbf{x}_i, \hat{\mathbf{w}})}_{\text{surprise of } y_i} \right)}^{\text{DL of observations given model}} + \frac{d}{2} \log_2(n)$$

- The sum of surprises of all observations is the **description length** of the observations given the (most probable) model in \mathcal{F} .

Intuitive interpretation of BIC

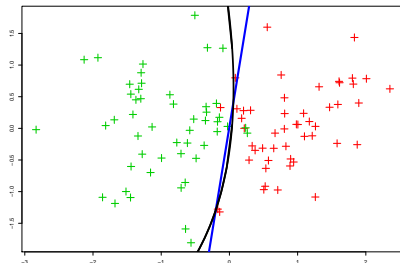
- Second term: description length of the model. Intuitive explanation:
 - ▶ The model, i.e. $\hat{\mathbf{w}} \in \mathbb{R}^d$, was estimated based on n samples.
 - ▶ Can quantize every component into \sqrt{n} levels. Why?
 - ▶ Remember the **standard parametric rate**:
 $1/\sqrt{n}$ represents the magnitude of the estimation error
 \rightsquigarrow **no need for encoding with greater precision.**
 - ▶ Grid of $(\sqrt{n})^d$ possible values for describing a model.
 - ▶ We need $\log_2((\sqrt{n})^d) = \log_2 n^{(d/2)} = (d/2) \log_2 n$ bits to encode $\hat{\mathbf{w}}$.
- In summary: $-\text{BIC} = \text{DL}(\text{data}|\text{model}) + \text{DL}(\text{model})$.
- Maximizing BIC = minimizing joint DL of data and model
 \rightsquigarrow **Minimum Description Length principle.**

Example: Bayesian logistic regression

Example: polynomial logistic regression, $n = 100$.

$\phi_1(\mathbf{x}) = (1, x_1, x_2)^t$, $\phi_2(\mathbf{x}) = (1, x_1, x_2, (x_1 + x_2)^2)^t$.

$$-\text{BIC} = \sum_{i=1}^n (-\log_2 p(y_i | \mathbf{x}_i, \hat{\mathbf{w}})) + \frac{d}{2} \log_2(n)$$



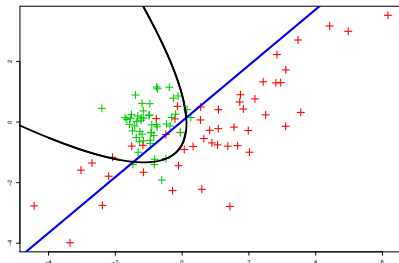
degree	#(param)	DL(data)	DL(model)	BIC score
1	3	16.36 bits	9.97 bits	-26.33
2	4	15.77 bits	13.29 bits	-29.06

Example: Bayesian logistic regression

Example: polynomial logistic regression, $n = 100$.

$$\phi_1(\mathbf{x}) = (1, x_1, x_2)^t, \quad \phi_2(\mathbf{x}) = (1, x_1, x_2, (x_1 + x_2)^2)^t.$$

$$-\text{BIC} = \sum_{i=1}^n (-\log_2 p(y_i | \mathbf{x}_i, \hat{\mathbf{w}})) + \frac{d}{2} \log_2(n)$$



degree	#(param)	DL(data)	DL(model)	BIC score
1	3	58.56 bits	9.97 bits	-68.53
2	4	38.05 bits	13.29 bits	-51.34

Subsection 3

Sparse models

Sparse Models

- Sometimes, we have many more dimensions d than training cases n .
- Corresponding design matrix X is “short and fat”, rather than “tall and skinny”.
- This is called **small n , large d problem**.
- For example, with **gene microarrays**, it is common to measure the expression levels of $d \approx 20,000$ genes, but to only get $n \approx 100$ samples (for instance, from 100 patients).
- Q: what is the **smallest set of features** that can accurately predict the response in order to **prevent overfitting**, to **reduce the cost** of building a diagnostic device, or to help with **scientific insight** into the problem?

Bayesian variable selection

- Let $\gamma_j = 1$ if feature j is **relevant**, and let $\gamma_j = 0$ otherwise.
- Our goal is to compute the posterior over models

$$p(\gamma|\mathcal{D}) = \frac{\exp(-f(\gamma))}{\sum_{\gamma'} \exp(-f(\gamma'))},$$

where $f(\gamma)$ is the cost function:

$$f(\gamma) = -[\log p(\mathcal{D}|\gamma) + \log p(\gamma)].$$

- For example, suppose we generate $n = 20$ samples from a $d = 10$ dimensional linear regression model, $y_i \sim N(w^t x_i, \sigma^2)$, in which $K = 5$ elements of w are non-zero.
- Enumerate all $2^{10} = 1024$ models and compute $p(\gamma|\mathcal{D})$ for each one.

Bayesian variable selection

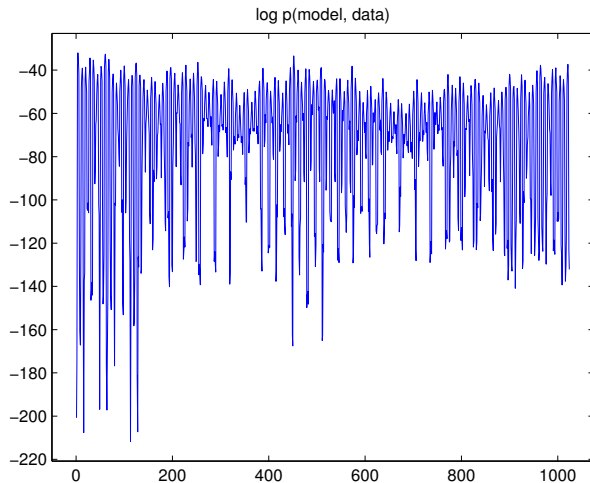


Fig 13.1 in K. Murphy: Score function $f(\gamma)$ for all possible models.

Bayesian variable selection

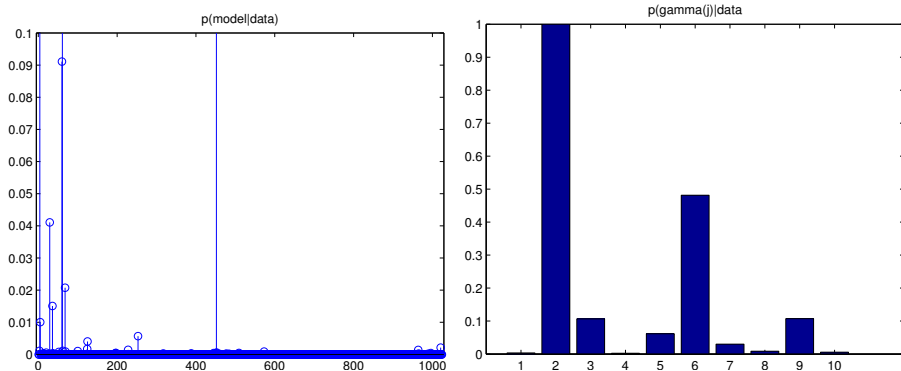


Fig 13.1 in K. Murphy. Left: Posterior over all 1024 models. Vertical scale has been truncated at 0.1 for clarity. Right: Marginal inclusion probabilities $p(\gamma_j = 1|\mathcal{D})$. The true model is $\{2, 3, 6, 8, 9\}$

Bayesian variable selection

- Interpreting the posterior over a large number of models is difficult
 \rightsquigarrow seek **summary statistics**.

- A natural one is the **posterior mode**, or MAP estimate

$$\hat{\gamma} = \arg \max p(\gamma|\mathcal{D}) = \arg \min f(\gamma).$$

- However, the mode is often not representative of the full posterior mass. A better summary is the **median model**, computed using

$$\hat{\gamma} = \{j : p(\gamma_j = 1|\mathcal{D}) > 0.5\}$$

This requires computing the **posterior marginal inclusion probabilities** $p(\gamma_j = 1|\mathcal{D})$.

Bayesian variable selection

- The above example illustrates the **gold standard** for variable selection: the problem was small ($d = 10$)
 \rightsquigarrow we were able to **compute the full posterior exactly**.
- Of course, variable selection is most useful in the cases where the number of dimensions is **large**.
- There are 2^d possible models (bit vectors) \rightsquigarrow **impossible** to compute the full posterior in general.
- Even finding summaries (MAP, or marginal inclusion probabilities) is **intractable**
 \rightsquigarrow **algorithmic speedups** necessary.
- But first, focus on the computation of $p(\gamma|\mathcal{D})$.

The spike and slab model

- The posterior is given by

$$p(\gamma|\mathcal{D}) \propto p(\gamma)p(\mathcal{D}|\gamma)$$

- It is common to use the following prior:

$$p(\gamma) = \prod_{j=1}^d \text{Ber}(\gamma_j|\pi_0) = \pi_0^{\|\gamma\|_0} (1 - \pi_0)^{d - \|\gamma\|_0},$$

$$\log p(\gamma|\pi_0) = -\lambda \|\gamma\|_0 + \text{const.},$$

where π_0 is the probability that a feature is relevant,

and $\|\gamma\|_0 = \sum_{j=1}^d \gamma_j$ is the ℓ_0 pseudo-norm, i.e., the **number of non-zero elements**.

- $\lambda = \log \frac{1-\pi_0}{\pi_0}$ controls the **sparsity** of the model.
- Setting $\sigma^2 = 1$, we can write the likelihood as follows:

$$p(\mathcal{D}|\gamma) = p(\mathbf{y}|X, \gamma) = \int p(\mathbf{y}|X, \mathbf{w}, \gamma) p(\mathbf{w}|\gamma) d\mathbf{w}$$

The spike and slab model

- Prior $p(\mathbf{w}|\gamma)$. If $\gamma_j = 0$, feature j is **irrelevant**, so we expect $w_j = 0$. If $\gamma_j = 1$, we expect w_j to be non-zero.
- Standardized inputs \rightsquigarrow reasonable **prior** is $N(0, \sigma_w^2)$, where σ_w^2 reflects our expectation of the coefficients associated with the **relevant variables**:

$$p(w_j|\gamma_j) = \begin{cases} \delta_0(w_j) & , \text{ if } \gamma_j = 0 \\ N(w_j|0, \sigma_w^2) & , \text{ else} \end{cases}$$

- The first term is a **spike** at the origin.
- As $\sigma_w^2 \rightarrow \infty$, the distribution $p(w_j|\gamma_j = 1)$ approaches a **uniform** distribution \rightsquigarrow **slab** of constant height.
- **Spike and slab model** (Mitchell and Beauchamp 1988).
- Full Bayesian treatment is computationally challenging!

Simplifying the model

- Assume $\sigma_w^2 \rightarrow \infty$ (\rightsquigarrow uniform prior $p(w_j|\gamma_j)$ over nonzero components) and approximate the likelihood using **BIC**:

$$\begin{aligned}\log p(\mathcal{D}|\gamma) &= \int p(\mathbf{y}|X, \mathbf{w}, \gamma) p(\mathbf{w}|\gamma) d\mathbf{w} \\ &\approx \log p(\mathbf{y}|X, \hat{\mathbf{w}}_\gamma) - \frac{1}{2} \underbrace{\|\hat{\mathbf{w}}_\gamma\|_0}_{\text{degrees of freedom}} \log n,\end{aligned}$$

where $\hat{\mathbf{w}}_\gamma$ is the ML estimate.

- Another view of this model: minimize the negative log likelihood under a ℓ_0 constraint (or penalty in the Lagrangian form)

$$\text{minimize } -\log p(\mathbf{y}|X, \mathbf{w}) + \lambda \|\mathbf{w}\|_0.$$

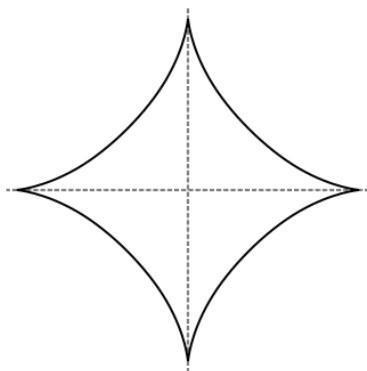
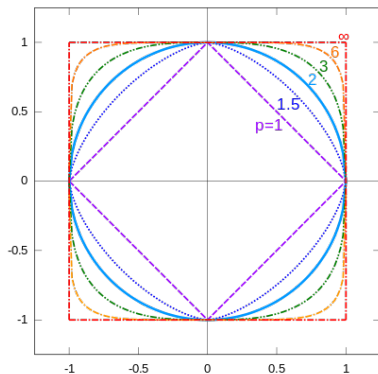
- Practical problem: ℓ_0 is highly non-convex!

Vector norms

The **vector p -norms** (**ℓ_p norms**) are defined by

$$\|\mathbf{x}\|_p = \left(\sum_{i=1}^n |x_i|^p \right)^{1/p}, \quad 1 \leq p \leq \infty,$$

$$\|\mathbf{x}\|_\infty = \max(|x_1|, \dots, |x_n|).$$



Simplifying the model further

- When we have many variables, it is computationally difficult to find the posterior mode
- Idea: replace discrete variables with continuous ones. Use continuous priors that “encourage” $w_j = 0$ by putting a lot of probability density near the origin, such as a zero-mean Laplace distribution.

$$p(\mathbf{w}|\lambda) = \prod_{j=1}^d \text{Lap}(w_j|0, 1/\lambda) \propto \prod_{j=1}^d \exp(-\lambda|w_j|)$$

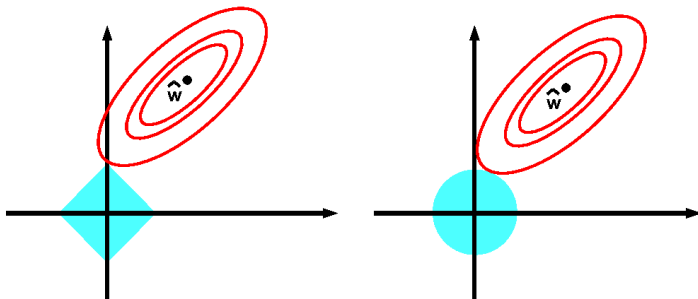
- Let us perform MAP estimation with this prior:

$$f(\mathbf{w}) = -\log p(\mathcal{D}|\mathbf{w}) - \log p(\mathbf{w}|\lambda) = NLL(\mathbf{w}) + \lambda\|\mathbf{w}\|_1.$$

where $\|\mathbf{w}\|_1 = \sum_{j=1}^d |w_j|$ is the ℓ_1 norm of \mathbf{w} and *NLL* means **negative log-likelihood**.

The Lasso

- For suitably large λ , the estimate $\hat{\mathbf{w}}$ will be sparse.
- Can be thought of as a convex approximation to the non-convex ℓ_0 objective.
- This model has the colorful name **least absolute shrinkage and selection operator**.
- For linear regression, $NLL(\mathbf{w}) = RSS(\mathbf{w})$,
a.k.a. **basis pursuit denoising** (Chen et al. 1998).



The Lasso

- Unfortunately, the $\|\mathbf{w}\|_1$ term is not differentiable at 0
 \rightsquigarrow non-smooth optimization problem.
- The **subderivative or subgradient** of a (convex) function $f : \mathcal{I} \rightarrow \mathbb{R}$ at a point x_0 is a scalar c such that

$$f(x) - f(x_0) \geq c(x - x_0), \quad \forall x \in \mathcal{I}$$

where \mathcal{I} is some interval containing x_0 .

Note that c is a **linear lower bound** to f at x_0 .

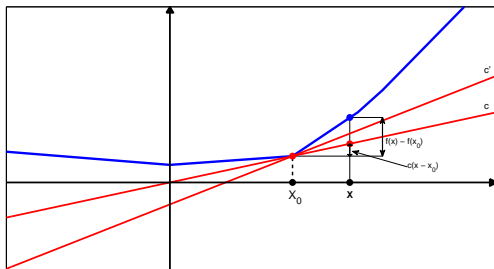


Fig. 13.4 in (K. Murphy)

The Lasso

- The set of all subderivatives is called the **subdifferential**
- For the **absolute value function** $f(x) = |x|$:

$$\partial f(x) = \begin{cases} -1 & , \text{ if } x < 0 \\ [-1, 1] & , \text{ if } x = 0 \\ +1 & , \text{ if } x > 0 \end{cases}$$

- For least-squares regression, it is easy to show that

$$\frac{\partial}{\partial w_j} \text{RSS}(\mathbf{w}) = a_j w_j - c_j$$

$$a_j = 2 \sum_{i=1}^n x_{ij}^2$$

$$c_j = 2 \sum_{i=1}^n x_{ij} (y_i - \mathbf{w}_{-j}^t \mathbf{x}_{i,-j}).$$

where \mathbf{w}_{-j} is \mathbf{w} without component j .

The Lasso

- c_j is (proportional to) the correlation between the j 'th feature \mathbf{x}_j and the residual due to other features, $r_{-j} = y - \mathbf{x}_{-j}^t \mathbf{w}_{-j}$.
- The magnitude of c_j is an indication of **how relevant** feature j is for predicting y .
- Adding the ℓ_1 penalty term:

$$\begin{aligned}\partial_{w_j} f(\mathbf{w}) &= (a_j w_j - c_j) + \lambda \partial_{w_j} \|\mathbf{w}\|_1 \\ &= \begin{cases} a_j w_j - c_j - \lambda & , \text{ if } w_j < 0 \\ [-c_j - \lambda, -c_j + \lambda] & , \text{ if } w_j = 0 \\ a_j w_j - c_j + \lambda & , \text{ if } w_j > 0 \end{cases}\end{aligned}$$

The Lasso

- Depending on the value of c_j , the solution to $\partial_{w_j} f(\mathbf{w}) = 0$ can occur at 3 different values of w_j :

$$\hat{w}_j = \begin{cases} (c_j + \lambda)/a_j & , \text{ if } c_j < -\lambda \\ 0 & , \text{ if } c_j \in [-\lambda, \lambda] \\ (c_j - \lambda)/a_j & , \text{ if } c_j > \lambda \end{cases}$$

- We can write this as follows:

$$\hat{w}_j = \text{soft}\left(\frac{c_j}{a_j}; \frac{\lambda}{a_j}\right),$$

where $\text{soft}(a; \delta) = \text{sign}(a)(|a| - \delta)_+$ and $x_+ = \max(x, 0)$ is the positive part of x .

- This is called **soft thresholding**.

The Lasso

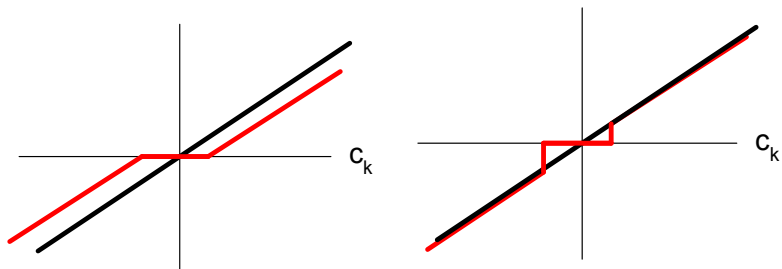


Fig. 13.5 in (K. Murphy). Black line: **Least squares fit** $w_k = c_k/a_k$. The red line (the regularized estimate) $\hat{w}_k(c_k)$, shifts the black line down (or up) by λ , except when $-\lambda \leq c_k \leq \lambda$, in which case it sets $w_k = 0$. By contrast, **hard thresholding** sets values of w_k to 0 if $-\lambda \leq c_k \leq \lambda$, but it **does not shrink the values of w_k outside of this interval**.

Lasso Algorithms: Coordinate-wise Descent

Sometimes it is hard to optimize all variables simultaneously, but it is easy to optimize them one by one.

Can solve for j -th coefficient w_j with all other coefficients held fixed:

$$\hat{w}_j = \arg \min_z f(\mathbf{w} + z \mathbf{e}_j),$$

where \mathbf{e}_j is the j -th unit vector. Cycle (potentially many times) through these component-wise updates:

for $j = 1, \dots, d$ **do**:

$$a_j = 2 \sum_{i=1}^n x_{ij}^2$$

$$c_j = 2 \sum_{i=1}^n x_{ij} (y_i - \mathbf{w}_{-j}^t \mathbf{x}_{i,-j})$$

$$w_j = \text{soft} \left(\frac{c_j}{a_j}; \frac{\lambda}{a_j} \right).$$