Chapter 2

Least squares problems



Linear curve fitting

- Notation: n objects at locations $x_i \in \mathbb{R}^p$. Every object has measurement $y_i \in \mathbb{R}$.
- Approximate "regression targets" y as a parametrized function of x.
- Consider a 1-dim problem initially.
- Start with n data points $(x_i, y_i), i = 1, \ldots, n$.
- Choose d basis functions $g_0(x), g_1(x), \ldots$
- Fitting to a line uses two basis functions $g_0(x) = 1$ and $g_1(x) = x$. In most cases $n \gg d$.
- Fit function = linear combination of basis functions: $f(x; w) = \sum_{j} w_{j}g_{j}(x) = w_{0} + w_{1}x.$
- $f(x_i) = y_i$ exactly is (usually) **not possible**, so approximate $f(x_i) \approx y_i$
- *n* residuals are defined by $r_i = y_i f(x_i) = y_i (w_0 + w_1 x_i)$.



Calculus or algebra?

- Quality of fit can be measured by residual sum of squares $RSS = \sum_i r_i^2 = \sum_i [y_i - (w_0 + w_1 x_i)]^2.$
- Minimizing RSS with respect to w_1 and w_0 provides the **least-squares fit.**
- To solve the least squares problem we can
 - 1. set the derivative of RSS to zero
 - \rightsquigarrow calculus, or
 - 2. solve an over-determined system

 \rightsquigarrow algebra: $w_0 + w_1 x_i = y_i, i = 1, \dots, n$

- The results you get are...
 - mathematically the same, but
 - have different numerical properties.



Matrix-vector form

• Write $f(x) \approx y$ in matrix-vector form for n observed points as

$$\underbrace{\begin{bmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_n \end{bmatrix}}_{X} \underbrace{\begin{bmatrix} w_0 \\ w_1 \end{bmatrix}}_{\boldsymbol{w}} \approx \underbrace{\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}}_{\boldsymbol{y}}$$

• We minimize the sum of squared errors, which is the squared norm of the residual vector $\boldsymbol{r} = \boldsymbol{y} - X \boldsymbol{w}$:

$$RSS = \sum_{i=1}^{n} (y_i - (Xw)_i)^2 = \|y - Xw\|^2 = \|r\|^2 = r^t r.$$

• RSS = 0 only possible if all the data points lie on a line.

Basis functions

X has as many columns as there are basis functions. Examples:

- High-dimensional linear functions $x \in \mathbb{R}^p$, $g_0(x) = 1$ and $g_1(x) = x_1, g_2(x) = x_2, \dots, g_p(x) = x_p$. $X_{i\bullet} = g^t(x_i) = (1, -x_i^t -), \quad (i\text{-th row of } X)$ $f(x; w) = w^t g = w_0 + w_1 x_1 + \dots + w_p x_p.$
- **Document analysis:** Assume a fixed collection of words:

$$oldsymbol{x} = ext{text} ext{document}$$

 $g_0(oldsymbol{x}) = 1$
 $g_i(oldsymbol{x}) = \#(ext{occurrences} ext{of} i- ext{th} ext{word} ext{in} ext{document})$
 $f(oldsymbol{x};oldsymbol{w}) = oldsymbol{w}^t oldsymbol{g} = w_0 + \sum_{i \in ext{words}} w_i g_i(oldsymbol{x}).$

Solution by Calculus

$$RSS = \mathbf{r}^{t}\mathbf{r} = (\mathbf{y} - X\mathbf{w})^{t}(\mathbf{y} - X\mathbf{w})$$
$$= \mathbf{y}^{t}\mathbf{y} - \mathbf{y}^{t}X\mathbf{w} - \mathbf{w}^{t}X^{t}\mathbf{y} + \mathbf{w}^{t}X^{t}X\mathbf{w}$$
$$= \mathbf{y}^{t}\mathbf{y} - 2\mathbf{y}^{t}X\mathbf{w} + \mathbf{w}^{t}X^{t}X\mathbf{w}.$$

Minimization: set the gradient (vector of partial derivatives) to zero:

$$\nabla_{\boldsymbol{w}}RSS = \frac{\partial RSS}{\partial \boldsymbol{w}} \stackrel{!}{=} \boldsymbol{0}.$$

We need some properties of vector derivatives:

$$\partial (A \boldsymbol{x}) / \partial \boldsymbol{x} = A^t$$

 $\partial (\boldsymbol{x}^t A) / \partial \boldsymbol{x} = A$
 $\partial (\boldsymbol{x}^t A \boldsymbol{x}) / \partial \boldsymbol{x} = A \boldsymbol{x} + A^t \boldsymbol{x}$ (if A is square

Normal Equations

$$\frac{\partial RSS}{\partial \boldsymbol{w}} = \frac{\partial}{\partial \boldsymbol{w}} \left[\boldsymbol{y}^t \boldsymbol{y} - 2\boldsymbol{y}^t X \boldsymbol{w} + \boldsymbol{w}^t X^t X \boldsymbol{w} \right]$$
$$= -2X^t \boldsymbol{y} + \left[X^t X \boldsymbol{w} + (X^t X)^t \boldsymbol{w} \right]$$
$$= -2X^t \boldsymbol{y} + 2X^t X \boldsymbol{w} = \boldsymbol{0}$$

Normal equations: $X^{t}Xw = X^{t}y$.

Could solve this system. **But:** All solution methods based on normal equations are **inherently susceptible to roundoff errors**:

$$\begin{split} k(X) &= \sigma_{\max} / \sigma_{\min}, \text{ where } X^t X \boldsymbol{v}_i = \sigma_i^2 \boldsymbol{v}_i \\ k(X^t X) &= \mu_{\max} / \mu_{\min}, \text{ where } X^t X X^t X \boldsymbol{v}_i = \mu_i^2 \boldsymbol{v}_i \\ X^t X X^t X \boldsymbol{v}_i &= X^t X \sigma_i^2 \boldsymbol{v}_i = \sigma_i^4 \boldsymbol{v}_i \Rightarrow \mu_i = \sigma_i^2 \\ \Rightarrow k(X^t X) = k^2(X), \end{split}$$

The algebraic approach will avoid this problem!

From Calculus to Algebra

$$\frac{\partial RSS(\boldsymbol{w})}{\partial \boldsymbol{w}} = -2X^t \boldsymbol{y} + 2X^t X \boldsymbol{w} \stackrel{!}{=} \boldsymbol{0}$$
$$\Rightarrow X^t (\boldsymbol{y} - X\hat{\boldsymbol{w}}) = X^t \boldsymbol{r} = \boldsymbol{0} \quad \Rightarrow \boldsymbol{r} \in N(X^t).$$

- Every Xw is in column space C(X), residual r is in the orthogonal complement N(X^t) (left nullspace).
- Let \hat{y} be the orthogonal projection of y on C(X) $\rightsquigarrow y$ can be split into $\hat{y} \in C(X) + r \in N(X^t)$.



Algebraic interpretation

• $y = \hat{y} \in C(X) + r \in N(X^t) \rightsquigarrow$ Consider over-determined systems

$$Xm{w}=m{y}=\hat{m{y}}+m{r}$$
 (solution impossible, if $m{r}
eq m{0}$)

$$X\hat{\boldsymbol{w}} = \hat{\boldsymbol{y}}$$
 (solvable, since $\hat{\boldsymbol{y}} \in C(X)$!)

• The solution \hat{w} of $Xw = \hat{y}$ makes the error as small as possible:

$$\|X w - y\|^2 = \|X w - (\hat{y} + r)\|^2 = \|X w - \hat{y}\|^2 + \|r\|^2$$

Reduce $||X w - \hat{y}||^2$ to zero by solving $X \hat{w} = \hat{y}$ and choosing $w = \hat{w}$. Remaining error $||r||^2$ cannot be avoided, since $r \in N(X^t)$.

$$X^{t}X\hat{\boldsymbol{w}} = X^{t}\hat{\boldsymbol{y}} = X^{t}\boldsymbol{y} \quad \Rightarrow \quad \hat{\boldsymbol{w}} = (X^{t}X)^{-1}X^{t}\boldsymbol{y} \text{ (if } X^{t}X \text{ invertible).}$$

- The fitted values at the sample points are $\hat{y} = X\hat{w} = X(X^tX)^{-1}X^ty$.
- $H = X(X^tX)^{-1}X^t$ is called **hat matrix** (puts a "hat" on $\boldsymbol{y} \rightsquigarrow \hat{\boldsymbol{y}}$).

Algebraic interpretation

- Left nullspace $N(X^t)$ is orthogonal complement of column space C(X).
- H is orthogonal projection on C(X):

$$HX = X(X^{t}X)^{-1}X^{t}X = X, \quad HN(X^{t}) = 0.$$

• M = I - H is orthogonal projection on nullspace of X^t :

$$MX = (I - H)X = X - X = 0, \quad MN(X^{t}) = M.$$

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

The algebra of Least Squares: H creates fitted values: $\hat{y} = Hy \rightsquigarrow \hat{y} \in C(X)$ M creates residuals: $r = My \rightsquigarrow r \in N(X^t)$

Algebraic interpretation

$X^{t}X$ is invertible iff X has linearly independent columns.

Why? $X^t X$ has the same nullspace as X: (i) If $a \in N(X)$, then $Xa = \mathbf{0} \Rightarrow X^t Xa = \mathbf{0} \rightsquigarrow a \in N(X^t X)$. (ii) If $a \in N(X^t X)$, then $a^t X^t Xa = 0 \Leftrightarrow ||Xa||^2 = 0$, so Xa has length zero $\Rightarrow Xa = \mathbf{0}$. Thus, every vector in one nullspace is also in the other one.

So if $N(X) = \{0\}$, then $X^t X \in \mathbb{R}^{d \times d}$ has full rank d.

When X has independent columns, $X^{t}X$ is positive definite.

Why? $X^t X$ is clearly symmetric and invertible. To show: All eigenvalues > 0 $X^t X \boldsymbol{v} = \lambda \boldsymbol{v} \rightsquigarrow \boldsymbol{v}^t X^t X \boldsymbol{v} = \lambda \boldsymbol{v}^t \boldsymbol{v} \rightsquigarrow \lambda = \frac{\|X \boldsymbol{v}\|^2}{\|\boldsymbol{v}\|^2} > 0.$

SVD for Least-Squares

- Goal: Avoid numerical problems for normal equations: $X^{t}Xw = X^{t}y, \quad k(X^{t}X) = k^{2}(X).$
- Idea: Apply the **SVD** directly to $X_{n \times d}$.
- The squared norm of the residual is

$$RSS = \|\boldsymbol{r}\|^2 = \|X\boldsymbol{w} - \boldsymbol{y}\|^2$$
$$= \|USV^t\boldsymbol{w} - \boldsymbol{y}\|^2$$
$$= \|U(SV^t\boldsymbol{w} - U^t\boldsymbol{y})\|^2$$
$$= \|SV^t\boldsymbol{w} - U^t\boldsymbol{y}\|^2$$

Last equation: U is orthonormal $\rightsquigarrow ||U\boldsymbol{a}||^2 = \boldsymbol{a}^t U^t U \boldsymbol{a} = \boldsymbol{a}^t \boldsymbol{a} = ||\boldsymbol{a}||^2$.

• Minimizing RSS is equivalent to minimizing $\|S\boldsymbol{z} - \boldsymbol{c}\|^2$ where $\boldsymbol{z} = V^t \boldsymbol{w}$ and $\boldsymbol{c} = U^t \boldsymbol{y}$.

SVD and **LS**

Recall: Columns u_i of $U_{n \times n}$ with $\sigma_i > 0$ form a **basis of** C(X). Remaining columns form **basis of** $N(X^t)$:

$$\boldsymbol{c} = U^{t}\boldsymbol{y} = \underbrace{\begin{bmatrix} - & \boldsymbol{u}_{1}^{t} & - \\ - & \boldsymbol{u}_{2}^{t} & - \\ \vdots & \\ - & \boldsymbol{u}_{d}^{t} & - \\ 0 & 0 & 0 \\ \vdots & \\ 0 & 0 & 0 \end{bmatrix}}_{\begin{array}{c} \vdots & \\ y_{n-1} \\ y_{n} \end{array}} + \underbrace{\begin{bmatrix} 0 & 0 & 0 \\ \vdots & \\ 0 & 0 & 0 \\ - & \boldsymbol{u}_{d+1}^{t} & - \\ - & \boldsymbol{u}_{d+2}^{t} & - \\ \vdots & \\ - & \boldsymbol{u}_{n}^{t} & - \end{bmatrix}}_{\begin{array}{c} y_{1} \\ y_{2} \\ \vdots \\ \vdots \\ y_{n-1} \\ y_{n} \end{bmatrix}}$$

$$\begin{bmatrix} c_{1} \\ \vdots \\ c_{d} \\ 0 \\ \vdots \\ 0 \end{bmatrix} \in C(X) \quad \begin{bmatrix} 0 \\ \vdots \\ 0 \\ c_{d+1} \\ \vdots \\ c_{n} \end{bmatrix}} \in N(X^{t})$$

SVD and bases for the 4 subspaces



SVD and LS

• $\|\boldsymbol{r}\|^2 = \|S\boldsymbol{z} - \boldsymbol{c}\|^2$ written in blocks:

$$\left\| \begin{bmatrix} \sigma_{1} & 0 & \dots & 0 \\ 0 & \sigma_{2} & \dots & 0 \\ 0 & 0 & \dots & \sigma_{d} \\ \hline 0 & 0 & \dots & 0 \\ \vdots & \vdots & \vdots & \vdots \\ 0 & 0 & \dots & 0 \end{bmatrix} \begin{bmatrix} z_{1} \\ z_{2} \\ \vdots \\ z_{d} \end{bmatrix} - \begin{bmatrix} c_{1} \\ \vdots \\ c_{d} \\ c_{d+1} \\ \vdots \\ c_{n} \end{bmatrix} \right\|^{2}$$

• To choose z so that $||r||^2$ is minimal requires $z_i = c_i/\sigma_i, i = 1, ..., d$ $\rightsquigarrow r_1 = r_2 = \cdots = r_d = 0.$

- Unavoidable error: $RSS = \|r\|^2 = c_{d+1}^2 + c_{d+2}^2 + \dots + c_n^2$.
- For very small singular values, use zeroing. RSS will increase: One additional term (usually small): $RSS' = c_d^2 + c_{d+1}^2 + c_{d+2}^2 + \cdots + c_n^2$, but often significantly better precision (reduced condition number).

Classification

Classification: Find **class boundaries** based on training data $\{(x_1, y_1), \ldots, (x_n, y_n)\}$. Use boundaries to classify new items x^* . Here, y_i is a discrete class indicator (or "label"). Example: Fish-packing plant wants to automate the process of sorting fish on conveyor belt using optical sensing.



(Duda, Hart, Stork, 2001)



(Duda, Hart, Stork, 2001)

Linear Discriminant Analysis (Ronald Fisher, 1936)



Main Idea: Simplify the problem by projecting down to a 1-dim subspace. **Question:** How should we select the **projection vector**, which optimally discriminates between the different classes?

• Let m_j an estimate of the class means μ_j :

$$\boldsymbol{m}_y = \frac{1}{n_y} \sum_{\boldsymbol{x} \in \text{class } y} \boldsymbol{x}, \quad n_y = \#(\text{objects in class } y).$$

• Projected samples: $x'_i = w^t x_i$, i = 1, 2, ..., n. Projected means:

$$\tilde{m}_y = \frac{1}{n_y} \sum_{x \in \text{class } y} w^t x = w^t m_y.$$

• First part of separation criterion (two-class case):

$$\max_{\boldsymbol{w}} [\boldsymbol{w}^t(\boldsymbol{m}_1 - \boldsymbol{m}_2)]^2 = \max_{\boldsymbol{w}} [\tilde{m}_1 - \tilde{m}_2]^2.$$

There might still be considerable overlap...

 should also consider the scatter or variance.

Two Gaussians with the same mean distance, but different variances:



Excursion: The multivariate Gaussian distribution



Probability density function: $p(\boldsymbol{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{\sqrt{2\pi |\boldsymbol{\Sigma}|}} \exp(-\frac{1}{2}(\boldsymbol{x} - \boldsymbol{\mu})^t \boldsymbol{\Sigma}^{-1}(\boldsymbol{x} - \boldsymbol{\mu}))$

Excursion: The multivariate Gaussian distribution

Covariance

(also written "*co*-variance") is a measure of how much **two random variables vary together.** Can be positive, zero, or negative.



Sample covariance matrix $\widehat{\Sigma} = \frac{1}{n} \sum_{i=1}^{n} (x_i - \overline{x}) (x_i - \overline{x})^t$, with sample mean $\overline{x} = \frac{1}{n} \sum_{i=1}^{n} x_i = m$. If $m = \mathbf{0} \rightsquigarrow \widehat{\Sigma} = \frac{1}{n} X^t X$.

• Assume both classes are Gaussians with the same covariance matrix. Let Σ_W be an estimate of this "within class" covariance matrix:

$$\Sigma_y = \frac{1}{n_y} \sum_{\boldsymbol{x} \in \text{class } y} (\boldsymbol{x} - \boldsymbol{m}_y) (\boldsymbol{x} - \boldsymbol{m}_y)^t,$$

$$\Sigma_W = 0.5 (\Sigma_1 + \Sigma_2).$$

• Variance of projected data:

$$\begin{split} \tilde{\Sigma}_y &= \frac{1}{n_y} \sum_{\boldsymbol{x} \in \mathsf{class } y} (\boldsymbol{w}^t \boldsymbol{x} - \tilde{m}_y) (\boldsymbol{w}^t \boldsymbol{x} - \tilde{m}_y)^t \\ &= \frac{1}{n_y} \sum_{\boldsymbol{x} \in \mathsf{class } y} \boldsymbol{w}^t (\boldsymbol{x} - \boldsymbol{m}_y) (\boldsymbol{x} - \boldsymbol{m}_y)^t \boldsymbol{w} = \boldsymbol{w}^t \Sigma_y \boldsymbol{w} \\ \tilde{\Sigma}_W &= 0.5 (\tilde{\Sigma}_1 + \tilde{\Sigma}_2) = \boldsymbol{w}^t \Sigma_W \boldsymbol{w} \in \mathbb{R}_+ \end{split}$$

• Strategy: $\Delta_{\tilde{m}}^2 = (\tilde{m}_1 - \tilde{m}_2)^2$ should be large, $\tilde{\Sigma}_W$ small.

$$J(\boldsymbol{w}) = \frac{\Delta_{\tilde{m}}^2}{\tilde{\Sigma}_W} = \frac{\boldsymbol{w}^t (\boldsymbol{m}_1 - \boldsymbol{m}_2) (\boldsymbol{m}_1 - \boldsymbol{m}_2)^t \boldsymbol{w}}{\boldsymbol{w}^t \Sigma_W \boldsymbol{w}}.$$

$$\frac{\partial}{\partial \boldsymbol{w}} J(\boldsymbol{w}) = \frac{\partial}{\partial \boldsymbol{w}} \frac{\boldsymbol{w}^{t} \Sigma_{B} \boldsymbol{w}}{\boldsymbol{w}^{t} \Sigma_{W} \boldsymbol{w}} \stackrel{!}{=} 0$$

$$= -\frac{\boldsymbol{w}^{t} \Sigma_{B} \boldsymbol{w}}{(\boldsymbol{w}^{t} \Sigma_{W} \boldsymbol{w})^{2}} 2 \Sigma_{W} \boldsymbol{w} + \frac{1}{\boldsymbol{w}^{t} \Sigma_{W} \boldsymbol{w}} 2 \Sigma_{B} \boldsymbol{w}$$

$$\Rightarrow \frac{\boldsymbol{w}^{t} \Sigma_{B} \boldsymbol{w}}{\boldsymbol{w}^{t} \Sigma_{W} \boldsymbol{w}} (-\Sigma_{W} \boldsymbol{w}) + \Sigma_{B} \boldsymbol{w} = 0$$

$$\Rightarrow \Sigma_{B} \boldsymbol{w} = \frac{\boldsymbol{w}^{t} \Sigma_{B} \boldsymbol{w}}{\boldsymbol{w}^{t} \Sigma_{W} \boldsymbol{w}} \Sigma_{W} \boldsymbol{w} =: \lambda \Sigma_{W} \boldsymbol{w}$$

• Let Σ_W be non-singular:

$$\begin{bmatrix} \Sigma_W^{-1} & \underbrace{\Sigma_B} \end{bmatrix} \boldsymbol{w}_{\Delta_{\boldsymbol{m}} \Delta_{\boldsymbol{m}}^t \boldsymbol{w} \propto \Delta_{\boldsymbol{m}}} = \lambda \boldsymbol{w}, \quad \text{with} \quad \lambda = \frac{\boldsymbol{w}^t \Sigma_B \boldsymbol{w}}{\boldsymbol{w}^t \Sigma_W \boldsymbol{w}} = J(\boldsymbol{w}).$$

- Thus, w is an eigenvector of $\Sigma_W^{-1}\Sigma_B$, the associated eigenvalue is the objective function! Maximum: eigenvector with largest eigenvalue.
- Unscaled Solution: $\hat{\boldsymbol{w}} = \Sigma_W^{-1} \Delta_{\boldsymbol{m}} = \Sigma_W^{-1} (\boldsymbol{m}_1 \boldsymbol{m}_2).$
- This is the solution of the linear system $\Sigma_W w = m_1 m_2$.
- Σ_W is a covariance matrix \rightsquigarrow there is an underlying data matrix A such that $\Sigma_W \propto A^t A \rightsquigarrow$ potential numerical problems: squared condition number compared to A...

Discriminant analysis and least squares

Theorem: The LDA vector $\hat{w}^{LDA} = \Sigma_W^{-1}(m_1 - m_2)$ coincides with the solution of the LS problem $\hat{w}^{LS} = \arg \min_{\boldsymbol{w}} ||X\boldsymbol{w} - \boldsymbol{y}||^2$ if

$$\begin{split} n_1 &= \# \text{ samples in class } \mathbf{1} \\ n_2 &= \# \text{ samples in class } \mathbf{2} \\ X &= \begin{bmatrix} - & x_1^t & - \\ - & x_2^t & - \\ \vdots & \\ - & x_n^t & - \end{bmatrix}, \quad \mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}, \\ \text{with} \quad & \frac{1}{n} \sum_{i=1}^n x_i = \mathbf{m} = \mathbf{0} \quad (\text{i.e. origin in sample mean}), \\ y_i &= \begin{cases} +1/n_1, & \text{if } x_i \text{ in class } \mathbf{1} \\ -1/n_2, & \text{else.} \end{cases} \Rightarrow \sum_{i=1}^n y_i = \mathbf{0}. \end{split}$$

Discriminant analysis and least squares (cont'd)

- "Within" covariance $\Sigma_W \propto \sum_{m{x} \in {\sf class}\, y} (m{x} m{m}_y) (m{x} m{m}_y)^t$.
- "Between" covariance $\Sigma_B \propto ({m m}_1 {m m}_2) ({m m}_1 {m m}_2)^t$
- The sum of both is the "total covariance" $\Sigma_B + \Sigma_W = \Sigma_T$ $\Sigma_T \propto \sum_i x_i x_i^t = X^t X.$
- We know that ${m w}^{ extsf{LDA}} \propto \Sigma_W^{-1} ({m m}_1 {m m}_2) \rightsquigarrow \Sigma_W {m w}^{ extsf{LDA}} \propto ({m m}_1 {m m}_2).$
- Now $\Sigma_B w^{\text{LDA}} = (m_1 m_2)(m_1 m_2)^t w^{\text{LDA}} \rightsquigarrow \Sigma_B w^{\text{LDA}} \propto (m_1 m_2).$
- $\Sigma_T \boldsymbol{w}^{\mathsf{LDA}} = (\Sigma_B + \Sigma_W) \boldsymbol{w}^{\mathsf{LDA}} \rightsquigarrow \Sigma_T \boldsymbol{w}^{\mathsf{LDA}} \propto (\boldsymbol{m}_1 \boldsymbol{m}_2).$
- With $X^t X = \Sigma_T$, $X^t y = m_1 m_2$, we arrive at $w^{\text{LDA}} \propto \Sigma_T^{-1}(m_1 - m_2) = \Sigma_T^{-1} X^t y \propto (X^t X)^{-1} X^t y = w^{\text{LS}}$.

Chapter 2 Least squares problems

Application Example: Secondary Structure Prediction in Proteins







α-helix

By Thomas Shafee, https://commons.wikimedia.org/w/index.php?curid=52821069

Short historical Introduction

- Genetics as a natural science started in 1866: Gregor Mendel performed experiments that pointed to the existence of biological elements called genes.
- **Deoxy-ribonucleic acid (DNA)** isolated by **Friedrich Miescher** in 1869.
- 1944: Oswald Avery (and coworkers) identified DNA as the major carrier of genetic material, responsible for inheritance.
 Ribose: (simple) sugar molecule, deoxy-ribose → loss of oxygen atom.
 Nucleic acid: overall name for DNA and RNA (large biomolecules). Named for their initial discovery in nucleus of cells, and for presence of phosphate groups (related to phosphoric acid).



By Miranda19983 - Own work, CC BY-SA 4.0, https://commons.wikimedia.org/w/index.php?curid=84120486

Short historical Introduction

- 1953, Watson & Crick: **3-dimensional structure of DNA.** They inferred the method of **DNA replication.**
- 2001: first draft of the human genome published by the Human Genome Project and the company Celera.
- Many new developments, such as Next Generation Sequencing,
 Deep learning etc.



By RE73 - Own work, CC BY-SA 3.0, https://commons.wikimedia.org/w/index.php?curid=18862884

Base pairs and the DNA



By Madprime (talk \hat{A} contribs) - Own work, CC BY-SA 3.0,

https://commons.wikimedia.org/w/index.php?curid=1848174

- DNA composed of 4 basic molecules
 ~> nucleotides.
- Nucleotides are identical up to different **nitrogen base:** organic molecule with a nitrogen atom that has the chemical properties of a base (due to free electron pair at nitrogen atom).
- Each nucleotide contains **phosphate**, **sugar** (of deoxy-ribose type), and one of the 4 bases: **Adenine, Guanine, Cytosine, Thymine** (A,G,C,T).
- Hydrogen bonds between base pairs: $G \equiv C$, A = T.



By OpenStax - https://cnx.org/contents/FPtK1zmh@8.25:fEl3C8Ot@10/Preface, CC BY 4.0,

 $https://commons.wikimedia.org/w/index.php?curid{=}30131206$

The structure of DNA

- DNA molecule is **directional** due to asymmetrical structure of the sugars which constitute the skeleton: Each sugar is connected to the strand **upstream** in its 5th carbon and to the strand **downstream** in its 3rd carbon.
- DNA strand goes from 5' to 3'. The directions of the two complementary DNA strands are reversed to one another (→ **Reversed Complement**).



Adapted from https://commons.wikimedia.org/w/index.php?curid=30131206



By Zephyris - Own work, CC BY-SA 3.0, https://commons.wikimedia.org/w/index.php?curid=15027555

Replication of DNA

Biological process of producing two replicas of DNA from one original DNA molecule.

Cells have the distinctive property of division

→ DNA replication is most essential part for **biological inheritance**.

Unwinding \rightsquigarrow single bases exposed on each strand.

Pairing requirements are strict \rightsquigarrow single strands are templates for re-forming identical double helix (up to **mutations**).

DNA polymerase: enzyme that catalyzes the synthesis of new DNA.



Genes and Chromosomes

- In higher organisms, DNA molecules are packed in a chromosome.
- **Genome:** total genetic information stored in the chromosomes.
- Every cell contains a **complete set** of the genome, differences are due to variable **expression** of genes.
- A gene is a sequence of nucleotides that encodes the synthesis of a gene product.





https://commons.wikimedia.org/w/index.php?curid=20539140

Gene expression: Process of synthesizing a gene product (often a protein) → controls timing, location, and amount.



Transcription: making of an RNA molecule from DNA template. **Translation:** construction of amino acid sequence from RNA.

Almost no exceptions (\rightsquigarrow retroviruses)

 \Rightarrow

Transcription



By Kelvinsong - Own work, CC BY 3.0, https://commons.wikimedia.org/w/index.php?curid=23086203



https://commons.wikimedia.org/w/index.php?curid=9810855

Translation

- mRNA molecules are translated by **ribosomes**: Enzyme that links together amino acids.
- Message is read three bases at a time.
- Initiated by the first AUG codon (codon = nucleotide triplet).
- Covalent bonds (=sharing of electron pairs) are made between adjacent amino acids
 ⇒ growing chain of amino acids ("polypeptide").
- When a "stop" codon (UAA, UGA, UAG) is encountered, translation stops.





By Boumphreyfr - Own work, CC BY-SA 3.0, https://commons.wikimedia.org/w/index.php?curid=7200200

The genetic code

Standard genetic code											
1st base	2nd base										
		U	С			Α		base			
U	UUU		UCU		UAU		UGU		U		
	UUC	(Pne/F) Pnenylalanine	UCC	10	UAC	(Tyr/Y) Tyrosine	UGC	(Cys/C) Cysteine	С		
	UUA		UCA	(Ser/S) Serine	UAA ^(B)	Stop (Ochre)	UGA ^[B]	Stop (Opal)	A		
	UUG	(Leu/L) Leucine	UCG		UAG ^[8]	Stop (Amber)	UGG	(Trp/W) Tryptophan	G		
с	CUU		CCU	(Pro/P) Proline	CAU		CGU		U		
	CUC		ccc		CAC	(His/H) Histidine	CGC		с		
	CUA		CCA		CAA		CGA	(Arg/R) Arginine	A		
	CUG		CCG		CAG	(GIn/Q) Glutamine	CGG		G		
A	AUU	(IIe/I) Isoleucine	ACU	(Thr/T) Threonine	AAU		AGU		U		
	AUC		ACC		AAC	(Asn/N) Asparagine	AGC	(Ser/S) Serine	С		
	AUA		ACA		AAA		AGA		A		
	AUG ^[A]	(Met/M) Methionine	ACG		AAG	(Lys/K) Lysine	AGG	(Arg/R) Arginine	G		
G	GUU	(Val/V) Valine	GCU	(Ala/A) Alanine	GAU		GGU		U		
	GUC		GCC		GAC	(Asp/D) Aspartic acid	GGC		с		
	GUA		GCA		GAA		GGA	(Gly/G) Glycine	A		
	GUG		GCG		GAG	(Glu/E) Glutamic acid	GGG		G		

Wikipedia

Highly redundant: only 20 (or 21) amino acids formed from $4^3 = 64$ possible combinations.



By Dancojocari. https://commons.wikimedia.org/w/index.php?curid=9176441

Proteins

- Linear polymer of amino acids, linked together by peptide bonds. Average size ≈ 200 amino acids, can be over 1000.
- To a large extent, cells are made of proteins.
- Proteins determine shape and structure of a cell.
 Main instruments of molecular recognition and catalysis.
- **Complex structure** with four hierarchical levels.
 - 1. **Primary structure**: amino acid sequence.
 - 2. Different regions form locally regular secondary structures like α helices and β -sheets.
 - 3. **Tertiary structure**: packing such structures into one or several 3D *domains*.
 - 4. Several domains arranged in a quaternary structure.

Molecular recognition

Interaction between molecules through noncovalent bonding



Crystal structure of a short peptide L-Lys-D-Ala-D-Ala (bacterial cell wall precursor) bound to the antibiotic vancomycin through hydrogen bonds. By M stone, CC BY-SA 3.0, https://commons.wikimedia.org/w/index.php?curid=2327682

Catalysis

Increasing the rate of a chemical reaction by adding a substance \rightsquigarrow catalyst.



Reaction Progress

Wikipedia

Protein Structure: primary to quaternary



Durbin et al., Cambridge University Press

Structure is determined by the **primary sequence** and their **physicochemical interactions** in the medium. **Structure determines functionality.**

Secondary Structure

Secondary structure: two main types: β -sheet and α -helix



The School of Biomedical Sciences Wiki

Short range interactions in the AA chain are important for the secondary structure: α -helix performs a 100° turn per amino acid \rightsquigarrow full turn after 3.6 AAs. Formation of a helix mainly depends on interactions in a 4 AA window.

Example: Cytochrome C2 Precursor

Secondary structure (h=helix)

amino acid sequence

hhhhhhhhh MKKGFLAAGVFAAVAFASGAALAEGDAAAGEKVSKKCLACHTFDQGGANKVGPNLFGVFE hhhhhhhh hhhhhhhh NTAAHKDDYAYSESYTEMKAKGLTWTEANLAAYVKDPKAFVLEKSGDPKAKSKMTFKLTK

hhhhhhhhhhh

DDEIENVIAYLKTLK



Given: Examples of known helices and non-helices in several proteins \rightarrow training set

Goal: Predict, mathematically, the existence and position of α -helices in **new proteins.**

Classification of Secondary Structure

Idea: Use a **sliding window** to cut the AA chain into pieces. 4 AAs are enough to capture one full turn \rightsquigarrow choose window of size 5.

Decision Problem: Find function f(...) that predicts for each substring in a window the structure:

$$f(AADTG) = \begin{cases} "Yes", if the central AA belongs to an α -helix, "No", otherwise$$

Problem: How should we numerically encode a string like AADTG?

Simple encoding scheme: Count the number of occurrences of each **AA in the window.** First order approximation, neglects AA's position within the window.

Example

- ...RAADTGGSDP...
- ...xxxhhhhhhx...
- $\dots x x x h h h h h h x \dots$
- ...xxxhhhhhhx...

(black \doteq structure info about central AA; green \doteq know secondary structure; red \doteq sliding window)

Α	C	D		G		R	S	Т		Y	Label
2	0	1	0	0	0	1	0	1	0	0	"No"
2	0	1	0	1	0	0	0	1	0	0	"Yes"
1	0	1	0	2	0	0	0	1	0	0	"Yes"
:			:	:	:	:	:	:	:		:

This is a binary classification problem ~> use Linear Discriminant Analysis

Discriminant Analysis

Consider $X_{n \times d}$, with n = #(windows) and d = #(AAs) = 20(or 21), and the *n*-vector of class indicators y

$$X = \begin{bmatrix} 2 & 0 & 1 & \dots & 0 & \dots \\ 2 & 0 & 1 & \dots & 1 & \dots \\ 1 & 0 & 1 & \dots & 2 & \dots \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \end{bmatrix} = \begin{bmatrix} - & \boldsymbol{x}_1^t & - \\ - & \boldsymbol{x}_2^t & - \\ & \vdots & \\ - & \boldsymbol{x}_n^t & - \end{bmatrix}, \quad \boldsymbol{y} = \begin{bmatrix} " \operatorname{No"} \\ " \operatorname{Yes"} \\ " \operatorname{Yes"} \\ \vdots \end{bmatrix}$$

For the binary class idicators, we use some numerical encoding scheme.

Interpretation with basis functions:

 $oldsymbol{x} =$ sequence of characters from alphabet \mathcal{A} $g_i(oldsymbol{x}) = \#($ occurences of letter i in sequence) $f(oldsymbol{x};oldsymbol{w}) = oldsymbol{w}^t oldsymbol{g} = \sum_{i \in \text{characters}} w_i g_i(oldsymbol{x})$

Discriminant analysis and least squares

Recall: The LDA vector $\hat{w}^{LDA} = \Sigma_W^{-1}(m_1 - m_2)$ coincides with the solution of the LS problem $\hat{w}^{LS} = \arg \min_{\boldsymbol{w}} ||X\boldsymbol{w} - \boldsymbol{y}||^2$ if

$$n_{1} = \# \text{ samples in class } \mathbf{1}$$

$$n_{2} = \# \text{ samples in class } \mathbf{2}$$

$$X = \begin{bmatrix} - & x_{1}^{t} & - \\ - & x_{2}^{t} & - \\ & \vdots \\ - & x_{n}^{t} & - \end{bmatrix}, \quad \mathbf{y} = \begin{bmatrix} y_{1} \\ y_{2} \\ \vdots \\ y_{n} \end{bmatrix},$$
ith
$$\frac{1}{n} \sum_{i=1}^{n} \mathbf{x}_{i} = \mathbf{m} = \mathbf{0} \text{ (i.e. origin in sample mean),}$$

$$y_{i} = \begin{cases} +1/n_{1}, & \text{if } \mathbf{x}_{i} \text{ in class } \mathbf{1} \\ -1/n_{2}, & \text{else.} \end{cases} \Rightarrow \sum_{i=1}^{n} y_{i} = 0$$

W

Singular Value Decomposition (SVD)

Recall: SVD for nonsquare matrix $X \in \mathbb{R}^{n \times d}$: $X = USV^t$.

Residual sum of squares: $RSS = \|\boldsymbol{r}\|^2 = \|X\boldsymbol{w} - \boldsymbol{y}\|^2 = \|USV^t\boldsymbol{w} - \boldsymbol{y}\|^2 = \|S\underbrace{V^t\boldsymbol{w}}_{\boldsymbol{z}} - \underbrace{U^t\boldsymbol{y}}_{\boldsymbol{c}}\|^2$

Minimizing $\|\boldsymbol{r}\|^2$ is equivalent to minimizing $\|S\boldsymbol{z} - \boldsymbol{c}\|^2$:

minimize
$$\|r\|^2 = \left\| \begin{bmatrix} \sigma_1 & 0 \\ \vdots & \ddots \\ 0 & \sigma_d \\ 0 & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & 0 \end{bmatrix} \cdot \begin{bmatrix} z_1 \\ \vdots \\ z_d \end{bmatrix} - \begin{bmatrix} c_1 \\ \vdots \\ c_d \\ c_{d+1} \\ \vdots \\ c_n \end{bmatrix} \right\|^2$$

We now choose z_k so that $\|\boldsymbol{r}\|^2$ is minimal, i.e., for $\sigma_k > 0$:

$$z_k = \frac{c_k}{\sigma_k}$$

Iterative Algorithm

In our problem we have d = 20 (or 21) and n > 10000. **Goal**: Use only $X^t X \in \mathbb{R}^{d \times d}$ and $X^t y \in \mathbb{R}^d$. **Initialize** $X^t X = 0$ (zero matrix), $X^t y = 0$. **Update:** for j = 1 to n:

$$X^{t}X + \boldsymbol{x}_{j}\boldsymbol{x}_{j}^{t} \longrightarrow X^{t}X$$
$$X^{t}\boldsymbol{y} + \boldsymbol{x}_{j}y_{j} \longrightarrow X^{t}\boldsymbol{y}$$

The first update procedure is correct, since

$$(X^{t}X)_{ik} = \sum_{j=1}^{n} x_{ji}x_{jk}$$

$$\Rightarrow X^{t}X = \sum_{j=1}^{n} \begin{bmatrix} x_{j1} \\ x_{j2} \\ \vdots \\ x_{jd} \end{bmatrix} \cdot [x_{j1}, x_{j2}, \dots, x_{jd}] = \sum_{j=1}^{n} x_{j}x_{j}^{t}$$

Iterative Algorithm

A similar calculation yields the other equation:

$$(X^{t}\boldsymbol{y})_{i} = \sum_{j} x_{ji}y_{j} \Rightarrow X^{t}\boldsymbol{y} = \sum_{j} \begin{bmatrix} x_{j1} \\ x_{j2} \\ \vdots \\ x_{jd} \end{bmatrix} \cdot y_{j} = \sum_{j=1}^{n} \boldsymbol{x}_{j}y_{j}$$

One remaining problem: In LDA we assumend that X was centered, i.e. the column sums are all zero. Compute the column sums as:

$$\mathbf{1}^{t}X = [1, 1, \dots, 1] \begin{bmatrix} - & \mathbf{x}_{1}^{t} & - \\ - & \mathbf{x}_{2}^{t} & - \\ & \vdots & \\ - & \mathbf{x}_{n}^{t} & - \end{bmatrix} = n \cdot [m_{1}, m_{2}, \dots, m_{d}] = n \cdot \mathbf{m}^{t}$$

 \rightsquigarrow "centered" $X_c = X - \mathbf{1}m^t = X - \frac{1}{n}\mathbf{1}\mathbf{1}^t X$

Centering

$$X_{c} = X - \mathbf{1}\mathbf{m}^{t} = X - \frac{1}{n}\mathbf{1}\mathbf{1}^{t}X$$
$$X_{c}^{t}X_{c} = X^{t}X + \frac{1}{n^{2}}X^{t}\mathbf{1}\underbrace{\mathbf{1}^{t}\mathbf{1}}_{=n}\mathbf{1}^{t}X - \frac{1}{n}X^{t}\mathbf{1}\mathbf{1}^{t}X - \frac{1}{n}X^{t}\mathbf{1}\mathbf{1}^{t}X$$
$$= X^{t}X - \frac{1}{n}X^{t}\mathbf{1}\mathbf{1}^{t}X$$
$$= X^{t}X - n \cdot \mathbf{m}\mathbf{m}^{t}$$

Iteratively update the vector $n \cdot m$ for every x_i corresponding to a new window position: Initialize $n \cdot m = 0$ and update $n \cdot m \leftarrow n \cdot m + x_i$

What about $X^t y$? We should have used

$$X_c^t \boldsymbol{y} = (X - \boldsymbol{1}\boldsymbol{m}^t)^t \boldsymbol{y} = (X^t - \boldsymbol{m}\,\boldsymbol{1}^t)\boldsymbol{y} = X^t \boldsymbol{y} - \boldsymbol{m}\,\boldsymbol{1}^t \boldsymbol{y}$$

But by construction, y is orthogonal to $\mathbf{1} \rightsquigarrow \mathbf{1}^t y = 0$, so nothing needs to be done!

Iterative Algorithm

Goal: Solution which only requires $X_c^t X_c \in \mathbb{R}^{d \times d}$ and $X_c^t y \in \mathbb{R}^d$ alone (and does not use X_c or y explicitly).

We need:

- The matrix V (for computing $\hat{w} = Vz$) **Solution:** columns of V are the eigenvectors of $X_c^t X_c$, corresponding eigenvalues are λ_i , $i = 1, ..., n \Rightarrow \sigma_i^2 = \lambda_i$
- For the nonzero SVs, we need $z_i = (U^t y)_i / \sigma_i = \sigma_i (U^t y)_i / \sigma_i^2$ Solution:

$$X_c = USV^t \Rightarrow V^t X_c^t \boldsymbol{y} = V^t V S^t U^t \boldsymbol{y} = S^t U^t \boldsymbol{y}$$

$$\Rightarrow z_i = (U^t \boldsymbol{y})_i / \sigma_i = (V^t X_c^t \boldsymbol{y})_i / \sigma_i^2$$

So \boldsymbol{z} and finally $\hat{\boldsymbol{w}} = V\boldsymbol{z}$ can be computed from $X_c^t X_c$ and $X_c^t \boldsymbol{y}$ alone!

Chapter 2

Least squares problems

Least-squares and dimensionality reduction

Least-squares and dimensionality reduction

Given n data points in d dimensions:

$$X = \begin{bmatrix} - & \boldsymbol{x}_1^t & - \\ - & \boldsymbol{x}_2^t & - \\ - & \vdots & - \\ - & \boldsymbol{x}_n^t & - \end{bmatrix} \in \mathbb{R}^{n \times d}$$

Want to reduce dimensionality from d to k. Choose k directions w_1, \ldots, w_k , arrange them as columns in matrix W:

$$W = \begin{bmatrix} | & | & | \\ \boldsymbol{w}_1 & \boldsymbol{w}_2 & \dots & \boldsymbol{w}_k \\ | & | & | \end{bmatrix} \in \mathbb{R}^{d \times k}$$

Project $x \in \mathbb{R}^d$ down to $z = W^t x \in \mathbb{R}^k$. How to choose W?

Encoding-decoding model

The projection matrix W serves two functions:

- Encode: $\boldsymbol{z} = W^t \boldsymbol{x}, \ \boldsymbol{z} \in \mathbb{R}^k, \ z_j = \boldsymbol{w}_j^t \boldsymbol{x}.$
 - The vectors w_j form a basis of the projected space.
 - We will require that this basis is orthonormal, i.e. $W^tW = I$.

• Decode:
$$\tilde{x} = Wz = \sum_{j=1}^k z_j w_j, \ \tilde{x} \in \mathbb{R}^d.$$

- If k = d, the above orthonormality condition implies $W^t = W^{-1}$, and encoding can be undone without loss of information.
- If k < d, the decoded \tilde{x} can only approximate $x \rightarrow the$ reconstruction error will be nonzero.
- Note that we did not include an intercept term. Assumption: origin of coordinate system is in the sample mean, i.e. $\sum_i x_i = 0$.

Principal Component Analysis (PCA)

We want the reconstruction error $\|m{x} - ilde{m{x}}\|^2$ to be small.

Objective: minimize $\min_{W \in \mathbb{R}^{d \times k}: W^t W = I} \sum_{i=1}^n \| \boldsymbol{x}_i - W W^t \boldsymbol{x}_i \|^2$

Finding the principal components

Projection vectors are orthogonal \rightsquigarrow can treat them separately:

$$egin{aligned} & \min_{oldsymbol{w}:\, \|oldsymbol{w}\|=1} \sum_{i=1}^n \|oldsymbol{x}_i - oldsymbol{w} oldsymbol{w}^t x_i \|^2 \ & \sum_i \|oldsymbol{x}_i - oldsymbol{w} oldsymbol{w}^t x_i - 2 x_i^t oldsymbol{w} oldsymbol{w}^t x_i + x_i^t oldsymbol{w} oldsymbol{w}^t x_i] \ & = \sum_i [oldsymbol{x}_i^t x_i - oldsymbol{x}_i^t oldsymbol{w} oldsymbol{w}^t x_i x_i^t oldsymbol{w} \\ & = \sum_i x_i^t x_i - \sum_i oldsymbol{w}^t x_i x_i^t oldsymbol{w} \\ & = \sum_i x_i^t x_i - oldsymbol{w}^t (\sum_{i=1}^n x_i x_i^t) oldsymbol{w} \\ & = \sum_i x_i^t x_i - oldsymbol{w}^t X^t X oldsymbol{w}. \end{aligned}$$

Finding the principal components

- Want to maximize $\boldsymbol{w}^t X^t X \boldsymbol{w}$ under the constraint $\|\boldsymbol{w}\| = 1$
- Can also maximize the ratio $J(\boldsymbol{w}) = \frac{\boldsymbol{w}^t X^t X \boldsymbol{w}}{\boldsymbol{w}^t \boldsymbol{w}}$.
- Optimal projection w is the eigenvector of $X^t X$ with largest eigenvalue (compare handout on spectral matrix norm).
- We assumed ∑_i x_i = 0, i.e. the columns of X sum to zero.
 → compute SVD of "centered" matrix X_c
 → column vectors in W are eigenvectors of X^t_cX_c
 → they are the principal components.

Eigen-faces [Turk and Pentland, 1991]

- d = number of pixels
- Each $oldsymbol{x}_i \in \mathbb{R}^d$ is a face image
- x_{ij} = intensity of the *j*-th pixel in image *i*



Conceptual: We can lean something about the structure of face images. **Computational:** Can use z_i for efficient nearest-neighbor classification: Much faster when $k \ll d$.

Information retrieval: Latent Semantic Analysis [Deerwater, 1990]

- d = number of words in the vocabulary, say 10000.
- Each $\boldsymbol{x}_i \in \mathbb{R}^d$ is a vector of word counts
- $x_{ij} =$ frequency of word j in document i



How to measure similarity between two documents? Dot products $x_i^t x_j$ In such high-dimensional spaces most pairs of vectors are almost orthogonal \rightsquigarrow scalar products tend to be "noisy". If $k \ll d$, $z_i^t z_j$ is probably a better similarity measure than $x_i^t x_j$.