Machine Learning 2020

Volker Roth

Department of Mathematics & Computer Science University of Basel

28th April 2020

Volker Roth (University of Basel)

Machine Learning 2020

28th April 2020 1 / 28

- 4 回 ト - 4 回 ト

Section 7

Support Vector Machines and Kernels

3.5

< □ > <

Э

Structure on canonical hyperplanes

Theorem (Vapnik, 1982)

Let *R* be the radius of the smallest ball containing the points $\mathbf{x}_1, \ldots, \mathbf{x}_n$: $B_R(\mathbf{a}) = \{\mathbf{x} \in \mathbb{R}^d : \|\mathbf{x} - \mathbf{a}\| < R, \ \mathbf{a} \in \mathbb{R}^d\}$. The set of canonical hyperplane decision functions $f(\mathbf{w}, w_0) = sign\{\mathbf{w}^t \mathbf{x} + w_0\}$ satisfying $\|\mathbf{w}\| \le A$ has VC dimension h bounded by

 $h \le R^2 A^2 + 1.$

Intuitive interpretation: margin = 1/||w|| \rightarrow minimizing capacity(\mathcal{H}) corresponds to maximizing the margin.

$$R[f_n] \leq R_{\mathsf{emp}}[f_n] + \sqrt{rac{\mathsf{a}}{n}} \left(\mathsf{capacity}(\mathcal{H}) + \ln rac{b}{\delta}
ight)$$

→ Large margin classifiers.

<ロト <部ト <きト <きト = 目

SVMs

• When the training examples are **linearly separable** we can maximize the margin by minimizing the regularization term

$$\|\mathbf{w}\|^2/2 = \sum_{i=1}^d w_i^2/2$$

subject to the classification constraints

$$y_i[x_i^t w] - 1 \ge 0, \ i = 1, \dots, n.$$



• The solution is defined only on the basis of a subset of examples or **support vectors.**

SVMs: nonseparable case

 Modify optimization problem slightly by adding a penalty for violating the classification constraints:

n

minimize
$$\|\boldsymbol{w}\|^2/2 + C\sum_{i=1}^n \xi_i$$

subject to relaxed constraints

$$y_i[x_i^t w] - 1 + \xi_i \ge 0, \ i = 1, \dots, n.$$

• The
$$\xi_i \ge 0$$
 are called **slack variables**.



SVMs: nonseparable case

We can also write the SVM optimization problem more compactly as

$$C\sum_{i=1}^{n} \overbrace{(1-y_i[\boldsymbol{x}_i^t \boldsymbol{w}])^+}^{\xi_i} + \|\boldsymbol{w}\|^2/2,$$

where $(z)^+ = z$ if $z \ge 0$ and zero otherwise.

• This is equivalent to regularized empirical loss minimization

$$\underbrace{\frac{1}{n}\sum_{i=1}^{n}(1-y_i[\boldsymbol{x}_i^t\boldsymbol{w}])^+}_{R_{emp}}+\lambda\|\boldsymbol{w}\|^2,$$

where $\lambda = 1/(2nC)$ is the regularization parameter.

SVMs and LOGREG

• When viewed from the point of view of regularized empirical loss minimization, SVM and logistic regression appear quite similar:

SVM:
$$\frac{1}{n} \sum_{i=1}^{n} (1 - y_i [\boldsymbol{x}_i^t \boldsymbol{w}])^+ + \lambda \|\boldsymbol{w}\|^2$$

LOGREG:
$$\frac{1}{n} \sum_{i=1}^{n} -\log \overbrace{\sigma(y_i [\boldsymbol{x}_i^t \boldsymbol{w}])}^{P(y_i | \boldsymbol{x}_i, \boldsymbol{w})} + \lambda \|\boldsymbol{w}\|^2,$$

where $\sigma(z) = (1 + e^{-z})^{-1}$ is the logistic function.

SVMs and LOGREG

• The difference comes from how we penalize errors:

Both:
$$\frac{1}{n} \sum_{i=1}^{n} \operatorname{Loss}(y_i[\boldsymbol{x}_i^t \boldsymbol{w}]) + \lambda \|\boldsymbol{w}\|^2$$
,

• SVM: Loss
$$(z) = (1 - z)^+$$

• LOGREG: Loss $(z) = \log(1 + \exp(-z))$



SVMs: solution, Lagrange multipliers

• Back to the separable case: how do we solve

minimize_{**w**} $\|\mathbf{w}\|^2/2$ s.t. $y_i[\mathbf{x}_i^t\mathbf{w}] - 1 \ge 0, i = 1, \dots, n.$

• Represent the constraints as individual loss terms:

$$\sup_{\alpha_i \geq 0} \alpha_i (1 - y_i[\boldsymbol{x}_i^t \boldsymbol{w}]) = \begin{cases} 0, & \text{if } y_i[\boldsymbol{x}_i^t \boldsymbol{w}] - 1 \geq 0, \\ \infty, & \text{otherwise.} \end{cases}$$

• Rewrite the minimization problem:

$$\begin{aligned} \mininimize_{\boldsymbol{w}} & \|\boldsymbol{w}\|^2/2 + \sum_{i=1}^n \sup_{\alpha_i \ge 0} \alpha_i (1 - y_i[\boldsymbol{x}_i^t \boldsymbol{w}]) \\ = \mininimize_{\boldsymbol{w}} & \sup_{\alpha_i \ge 0} \left(\|\boldsymbol{w}\|^2/2 + \sum_{i=1}^n \alpha_i (1 - y_i[\boldsymbol{x}_i^t \boldsymbol{w}]) \right) \end{aligned}$$

SVMs: solution, Lagrange multipliers

 Swap maximization and minimization (technically this requires that the problem is convex and feasible ~>> Slater's condition):

$$\mininimize_{\boldsymbol{w}} \left[\sup_{\alpha_i \ge 0} \left(\|\boldsymbol{w}\|^2 / 2 + \sum_{i=1}^n \alpha_i (1 - y_i[\boldsymbol{x}_i^t \boldsymbol{w}]) \right) \right]$$
$$= \maxinize_{\alpha_i \ge 0} \left[\min_{\boldsymbol{w}} \left(\underbrace{\|\boldsymbol{w}\|^2 / 2 + \sum_{i=1}^n \alpha_i (1 - y_i[\boldsymbol{x}_i^t \boldsymbol{w}])}_{J(\boldsymbol{w};\alpha)} \right) \right]$$

 We have to minimize J(w; α) over parameters w for fixed Lagrange multipliers α_i ≥ 0.
 Simple, because J(w) is convex ↔ set derivative to zero
 → only one stationary point → global minimum.

SVMs: solution, Lagrange multipliers

• Find optimal **w** by setting the derivatives to zero:

$$\frac{\partial}{\partial \boldsymbol{w}} J(\boldsymbol{w}; \boldsymbol{\alpha}) = \boldsymbol{w} - \sum_{i} \alpha_{i} y_{i} \boldsymbol{x}_{i} = 0 \implies \hat{\boldsymbol{w}} = \sum_{i} \alpha_{i} y_{i} \boldsymbol{x}_{i}.$$

• Substitute the solution back into the objective and get (after some re-arrangements of terms):

$$\begin{aligned} \max_{\alpha_i \ge 0} \min_{\boldsymbol{w}} \left(\|\boldsymbol{w}\|^2 / 2 + \sum_{i=1}^n \alpha_i (1 - y_i[\boldsymbol{x}_i^t \boldsymbol{w}]) \right) \\ = \max_{\alpha_i \ge 0} \left(\|\hat{\boldsymbol{w}}\|^2 / 2 + \sum_{i=1}^n \alpha_i (1 - y_i[\boldsymbol{x}_i^t \hat{\boldsymbol{w}}]) \right) \\ = \max_{\alpha_i \ge 0} \left(\sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n y_i y_j \alpha_i \alpha_j \boldsymbol{x}_i^t \boldsymbol{x}_j \right) \end{aligned}$$

SVMs: summary

• Find optimal Lagrange multipliers $\hat{\alpha}_i$ by maximizing

$$\sum_{i=1}^{n} \alpha_{i} - \frac{1}{2} \sum_{i,j=1}^{n} y_{i} y_{j} \alpha_{i} \alpha_{j} \boldsymbol{x}_{i}^{t} \boldsymbol{x}_{j} \quad \text{subject to } \alpha_{i} \geq 0.$$

- Only $\hat{\alpha}_i$'s corresponding to **support vectors** will be non-zero.
- Make **predictions** on any new example **x** according to:

$$\operatorname{sign}(\boldsymbol{x}^t \hat{\boldsymbol{w}}) = \operatorname{sign}(\boldsymbol{x}^t \sum_{i=1}^n \hat{\alpha}_i y_i \boldsymbol{x}_i) = \operatorname{sign}(\sum_{i \in SV} \hat{\alpha}_i y_i \boldsymbol{x}^t \boldsymbol{x}_i).$$

- Observation: dependency on input vectors only via dot products.
- Later we will introduce the **kernel trick** for efficiently computing these dot products in implicitly defined feature spaces.

SVMs: formal derivation

• Convex optimization problem: an optimization problem

minimize
$$f(\mathbf{x})$$
(1)subject to $g_i(\mathbf{x}) \leq 0, \quad i = 1, \dots, m$ (2)

is convex if the functions $f, g_1 \dots g_m : \mathbb{R}^n \to \mathbb{R}$ are convex.

• The Lagrangian function for the problem is

$$\mathcal{L}(\mathbf{x},\lambda_0,...,\lambda_m) = \lambda_0 f(\mathbf{x}) + \lambda_1 g_1(\mathbf{x}) + ... + \lambda_m g_m(\mathbf{x}).$$

Karush-Kuhn-Tucker (KKT) conditions: For each point x̂ that minimizes f, there exist real numbers λ₀,..., λ_m, called Lagrange multipliers, that simultaneously satisfy:

1
$$\hat{\boldsymbol{x}}$$
 minimizes $\mathcal{L}(\boldsymbol{x}, \lambda_0, \lambda_1, \dots, \lambda_m)$,

- 2 $\lambda_0 \ge 0, \lambda_1 \ge 0, \dots, \lambda_m \ge 0$, with at least one $\lambda_k > 0$,
- 3 Complementary slackness: $g_i(\hat{\mathbf{x}}) < 0 \Rightarrow \lambda_i = 0, 1 \le i \le m$.

SVMs: formal derivation

- Slater's condition: If there exists a strictly feasible point z satisfying g₁(z) < 0, ..., g_m(z) < 0, then one can set λ₀ = 1.
- Assume that Slater's condition holds. Minimizing the supremum $\mathcal{L}^*(\mathbf{x}) = \sup_{\lambda \ge 0} \mathcal{L}(\mathbf{x}, \lambda)$, is the **primal problem P**:

$$\hat{\boldsymbol{x}} = \operatorname*{argmin}_{\boldsymbol{x}} \mathcal{L}^*(\boldsymbol{x}).$$

Note that

$$\mathcal{L}^*(\boldsymbol{x}) = \sup_{\lambda \ge 0} \left(f(\boldsymbol{x}) + \sum_{i=1}^m \lambda_i g_i(\boldsymbol{x}) \right) = \begin{cases} f(\boldsymbol{x}) & \text{, if } g_i(\boldsymbol{x}) \le 0 \,\forall i \\ \infty & \text{, else.} \end{cases}$$

 \rightsquigarrow Minimizing $\mathcal{L}^*(x)$ is equivalent to minimizing f(x).

• The maximizer of the **dual problem D** is

$$\hat{\lambda} = rgmax_{\lambda} \mathcal{L}_*(\lambda), ext{ where } \mathcal{L}_*(\lambda) = \inf_{m{x}} \mathcal{L}(m{x},\lambda).$$

・ロト ・ 一 ト ・ ヨ ト ・ ヨ ト ・ ヨ

SVMs: formal derivation

- The non-negative number min(P) max(D) is the **duality gap.**
- Convexity and Slater's condition imply strong duality:
 - The optimal solution (x̂, λ̂) is a saddle point of L(x, λ)
 The duality gap is zero.
- Discussion: For any real function f(a, b) min_a[max_b f(a, b)] ≥ max_b[min_a f(a, b)].
 Equality → saddle value exists.



By Nicoguaro - Own work, CC BY 3.0, https://commons.wikimedia.org/w/index.php?curid=20570051

・ロト ・ 一 ・ ・ ー ・

Kernel functions

- A **kernel function** is a real-valued function of two arguments, $k(\mathbf{x}, \mathbf{x}') \in \mathbb{R}$, for $\mathbf{x}, \mathbf{x}' \in \mathcal{X}$.
- Typically the function is **symmetric**, and sometimes non-negative.
- In the latter case, it might be interpreted as a measure of similarity.
- Example: isotropic Gaussian kernel:

$$k(\mathbf{x}, \mathbf{x}') = \exp\left(-\frac{\|\mathbf{x} - \mathbf{x}'\|^2}{2\sigma^2}\right)$$

Here, σ^2 is the bandwidth. This is an example of a radial basis function (RBF) kernel (only a function of $||\mathbf{x} - \mathbf{x}'||^2$).

Mercer kernels

• A symmetric kernel is a Mercer kernel, iff the Gram matrix

$$K = \begin{pmatrix} k(\mathbf{x}_1, \mathbf{x}_1) & \dots & k(\mathbf{x}_1, \mathbf{x}_n) \\ \vdots & \vdots \\ k(\mathbf{x}_n, \mathbf{x}_1) & \dots & k(\mathbf{x}_n, \mathbf{x}_n) \end{pmatrix}$$

- is **positive semidefinite** for any set of inputs $\{x_i, \ldots, x_n\}$.
- Mercer's theorem: Eigenvector decomposition

$$K = V \Lambda V^t = (V \Lambda^{1/2}) (V \Lambda^{1/2})^t =: \Phi \Phi^t.$$

Eigenvectors: columns of *V*. Eigenvalues: entries of diagonal matrix $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_n)$. Note that $\lambda_i \in \mathbb{R}$ and $\lambda_i \geq 0$. Define $\phi(\mathbf{x}_i)^t = i$ -th row of $\Phi = V_{[i \bullet]} \Lambda^{1/2}$ $\rightsquigarrow k(\mathbf{x}_i, \mathbf{x}_{i'}) = \phi(\mathbf{x}_i)^t \phi(\mathbf{x}_{i'})$.

• Entries of *K*: **inner product of some feature vectors**, implicitly defined by eigenvectors *V*.

・ロト ・ 一 ト ・ ヨ ト ・ ヨ ト ・ ヨ

Mercer kernels

• If the kernel is **Mercer**, then there exists $\phi : \mathbf{x} \to \mathbb{R}^d$ such that $k(\mathbf{x}, \mathbf{x}') = \phi(\mathbf{x})^t \phi(\mathbf{x}'),$

where ϕ depends on the eigenfunctions of k (d might be infinite).

• Example: Polynomial kernel

$$k(\mathbf{x},\mathbf{x}')=(1+\mathbf{x}^t\mathbf{x}')^m.$$

Corresponding feature vector contains terms up to degree m. Example: $m = 2, x \in \mathbb{R}^2$:

$$(1 + \mathbf{x}^{t}\mathbf{x}')^{2} = 1 + 2x_{1}x_{1}' + 2x_{2}x_{2}' + (x_{1}x_{1}')^{2} + (x_{2}x_{2}')^{2} + 2x_{1}x_{1}'x_{2}x_{2}'.$$

Thus,

$$\phi(\mathbf{x}) = [1, \sqrt{2}x_1, \sqrt{2}x_2, x_1^2, x_2^2, \sqrt{2}x_1x_2]^t.$$

Equivalent to working in a 6-dim feature space.

• Gaussian kernel: feature map lives in an infinite dimensional space.

< ロ > < 同 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ >

Kernels for documents

- In document classification or retrieval, we want to compare two documents, x_i and x_{i'}.
- Bag of words representation:
 x_{ij} is the number of times word *j* occurs in document *i*.
- One possible choice: **Cosine similarity:**

$$k(\mathbf{x}_i, \mathbf{x}_{i'}) = \frac{\mathbf{x}_i^t \mathbf{x}_{i'}}{\|\mathbf{x}_i\| \|\mathbf{x}_{i'}\|} =: \phi(\mathbf{x}_i)^t \phi(\mathbf{x}_{i'}).$$

Problems:

- Popular words (like "the" or "and") are not discriminative
 ~> remove these stop words.
- Bias: once a word is used in a document, it is very likely to be used again.
- Solution: Replace word counts with "normalized" representation.

Kernels for documents

• TF-IDF "term frequency inverse document frequency": **Term frequency** is log-transform of the count:

 $\mathsf{tf}(x_{ij}) = \mathsf{log}(1 + x_{ij})$

Inverse document frequency:

$$\operatorname{idf}(j) = \log rac{\#(\operatorname{documents})}{\#(\operatorname{documents containing term } j)} = \log rac{1}{\hat{p}_j}.$$

→→ Shannon information content:

idf is a measure of how much information a word provides

• Combine with tf \rightsquigarrow counts weighted by information content:

 $\mathsf{tf}\mathsf{-idf}(\mathbf{x}_i) = [\mathsf{tf}(\mathbf{x}_{ij}) \cdot \mathsf{idf}(j)]_{j=1}^V$, where $V = \mathsf{size}$ of vocabulary.

• We then use this inside the cosine similarity measure. With $\phi(\mathbf{x}) = \text{tf-idf}(\mathbf{x})$:

$$k(\mathbf{x}_i, \mathbf{x}_{i'}) = \frac{\phi(\mathbf{x}_i)^t \phi(\mathbf{x}_{i'})}{\|\phi(\mathbf{x}_i)\| \|\phi(\mathbf{x}_{i'})\|}.$$

String kernels

- Real power of kernels arises for structured input objects.
- Consider two strings x, and x' of lengths d, d', over alphabet A.
 Idea: define similarity as the number of common substrings.
- If s is a substring of $x \rightsquigarrow \phi_s(x) =$ number of times s appears in x.
- String kernel

$$k(x,x') = \sum_{s \in \mathcal{A}^*} w_s \phi_s(x) \phi_s(x'),$$

where $w_s \ge 0$ and $\mathcal{A}^* =$ set of all strings (any length) from \mathcal{A} .

- One can show: Mercer kernel, can be computed in O(|x| + |x'|) time using suffix trees (Shawe-Taylor and Cristianini, 2004).
- Special case: $w_s = 0$ for |s| > 1: **bag-of-characters kernel:** $\phi(x)$ is the number of times each character in \mathcal{A} occurs in x.

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 つくつ

The kernel trick

- Idea: modify algorithm so that it **replaces all inner products** $x^t x'$ with a call to the **kernel function** k(x, x').
- Kernelized ridge regression: $\hat{w} = (X^t X + \lambda I)^{-1} X^t y$. Matrix inversion lemma:

$$(I + UV)^{-1}U = U(I + VU)^{-1}$$

Define new variables α_i :

$$\hat{\boldsymbol{w}} = (X^t X + \lambda I)^{-1} X^t \boldsymbol{y}$$

= $X^t \underbrace{(XX^t + \lambda I)^{-1} \boldsymbol{y}}_{\hat{\alpha}} = \sum_{i=1}^n \hat{\alpha}_i \boldsymbol{x}_i.$

 \rightarrow solution is linear sum of the *n* training vectors.

The kernel trick

• Use this and the kernel trick to make predictions for x:

$$\hat{f}(\mathbf{x}) = \hat{\mathbf{w}}^t \mathbf{x} = \sum_{i=1}^n \hat{\alpha}_i \mathbf{x}_i^t \mathbf{x} = \sum_{i=1}^n \hat{\alpha}_i k(\mathbf{x}_i, \mathbf{x}).$$

Same for SVMs:

$$\hat{\boldsymbol{w}}^t \boldsymbol{x} = \sum_{i \in SV} \hat{\alpha}_i y_i \boldsymbol{x}_i^t \boldsymbol{x} = \sum_{i \in SV} \hat{\alpha}'_i k(\boldsymbol{x}_i, \boldsymbol{x})$$

• ...and for most other classical algorithms in ML!

Some applications in bioinformatics

• Bioinformatics: often non-vectorial data-types:



interaction graphs

- phylogenetic trees
- strings GSAQVKGHGKKVADALTNAVAHV
- Data fusion: convert data of each type into kernel matrix
 - \Rightarrow fuse kernel matrices
 - \Rightarrow "common language" for heterogeneous data.

HEMOGI OBIN

RBF kernels from expression data

- **Measurements** (for each gene): vector of expression values under different experimental conditions
- "classical" RBF kernel $k(x_1, x_2) = \exp(-\sigma ||x_1 x_2||^2)$



Diffusion kernels from interaction-graphs

- A: Adjacency matrix, D: node degrees, L = D A.
- $K := \frac{1}{Z(\beta)} \exp(-\beta L)$ with transition probabilities β .
- Physical interpretation (random walk): randomly choose next node among neighbors.
- Self-transition occurs with prob. $1 d_i \beta$







• K_{ij} : prob. for walk from *i* to *j*.

(Kondor and Lafferty, 2002)

< ロ > < 同 > < 三 > < 三 >

Alignment kernels from sequences

Alignment with Pair HMMs \rightsquigarrow Mercer kernel (Watkins, 2000). Image source: Durbin, Eddy, Krogh, Mitchison. Biological Seguence Alignment. Cambridge.



HBA_HUMAN GSAQVKGHGKKVADALTNAVAHV---D--DMPNALSALSDLHAHKL ++ ++++H+ KV + +A ++ +L+ L+++H+ K LGB2_LUPLU NNPELQAHAGKVFKLVYEAAIQLQVTGVVVTDATLKNLGSVHVSKG



Combination of heterogeneous data

Adding kernels \Rightarrow new kernel:

$$egin{aligned} &k_1(x,y)=\phi_1(x)\cdot\phi_1(y),\ &k_2(x,y)=\phi_2(x)\cdot\phi_2(y) \end{aligned} \Rightarrow &k'=k_1+k_2=(\phi_1(x))\cdot(\phi_1(y)),\ &\phi_2(y) \end{aligned}$$

Fusion & relevance determination: kernel-combinations

