

Linear Classification

**Perceptron Alg.** ! Work for linearly separable D ! Terminates after no more than  $1/\gamma^2$  updates, where  $\gamma = \min_{i=1,\dots,n} t_i \mathbf{w}_*^T \tilde{\phi}_i$  !  
Normalize features  $\phi_i \rightarrow \tilde{\phi}_i$  !  
Update rule (when misclassified  $t_i \mathbf{w}^T \tilde{\phi}_i \leq 0$ ):  $\mathbf{w} \leftarrow \mathbf{w} + t_i \tilde{\phi}_i$

Stochastic Gradient Descent

$\mathbf{w}_{k+1} = \mathbf{w}_k - \eta \nabla_{\mathbf{w}} E_i(\mathbf{w}_k)$

Least square estimation

$$E(\mathbf{w}) = \frac{1}{2} \sum_{i=1}^n (\mathbf{y}(\mathbf{x}_i) - t_i)^2$$
$$= \frac{1}{2} \|\Phi \mathbf{w} - \mathbf{t}\|^2$$

**Gradient for Squared Error** (yields normal equation for least squares if set to 0)

$$\nabla_{\mathbf{w}} E = \sum_{i=1}^n \frac{\partial E}{\partial y_i} \nabla_{\mathbf{w}} y_i = \Phi^T (\underbrace{\Phi \mathbf{w} - \mathbf{t}}_{residual})$$

Multi-Layer Perceptron

Forward pass

$$a_q^l = (\mathbf{w}_q^{(l)})^T \mathbf{z}^{(l-1)} + b_q^{(l)}$$
$$z_q^{(l)} = g(a_q^{(l)}), q = 1, \dots, h_l$$

Backward pass

$$r^{(L)} = \frac{\partial E_i}{\partial a^{(L)}} = \begin{cases} a^L - t_i & \text{for } E_{sq} \\ \sigma(a^L) - \hat{t}_i & \text{for } E_{log} \end{cases}$$

$$r_q^{(l)} = g'(a_q^{(l)}) \sum_{j=1}^{h_{l+1}} w_{jq}^{(l+1)} r_j^{(l+1)}$$

Gradient computation

$$\nabla_{w^{(l)}_i} E_i = r^{(l+1)} \mathbf{z}^{(l)}, \nabla_{w^{(l)}_q} E_i = r_q^{(1)} \mathbf{x}$$

For b: keep residual only ! (i.e  $\mathbf{x} = 1$ )

Momentum learning

$$\Delta \mathbf{w}_k = \mathbf{w}_{k+1} - \mathbf{w}_k$$

$$\Delta \mathbf{w}_k = -\eta(1 - \mu) \nabla_{\mathbf{w}_k} E + \mu \Delta \mathbf{w}_{k-1}$$

$\eta$  = learning rate e.g.  $1/k$ ,  $\mu$  = momentum term

Linear Regression. LSE

Univariate Linear Regression

$\langle x \rangle = n^{-1} \sum_i x_i$ . Similar for  $t, tx, x^2$

$$y_i = wx_i + b$$

$$\rightarrow w = \frac{Cov(x, t)}{Var(x)}, b = \langle t \rangle - w \langle x \rangle$$

Normal Equations  $(\Phi^T \Phi) \mathbf{w} = \Phi^T \mathbf{t}$

$$\hat{\mathbf{w}} = \underset{\mathbf{w}}{\operatorname{argmin}} E(\mathbf{w}) = (\Phi^T \Phi)^{-1} \Phi^T \mathbf{t}$$

Probability. Decision Theory

**Probability**  $E[X] = E[E[X|Y]]$   
Sum rule:  $P(X) = \sum_Y P(X, Y)$   
Product rule:  $P(X, Y) = P(X|Y)P(Y)$   
Bayes  $P(B|F) = \frac{P(F|B)P(B)}{P(F)}$   
 $Var(t) = E[Var(t|x)] + Var(E[t|x])$

Bayes-Optimal Classifier

$$f^*(\mathbf{x}) = \operatorname{argmax}_{t \in \tau} \underbrace{P(t|\mathbf{x})}_{p(\mathbf{x}|t)P(t)}$$

Bayes error:  $R = P\{f(x) \neq t\}$

$$R^* = R(f^*) = 1 - E \left[ \max_{k \in \tau} P(t = k|\mathbf{x}) \right]$$

Optimal discriminant:

$$y^*(x) = \log \frac{p(\mathbf{x}|t=1)}{p(\mathbf{x}|t=0)} + \log \frac{P(t=1)}{P(t=0)} > 0$$

Bayes under Loss function. Risk:

$R(f) = E[L(f(\mathbf{x}), t)]$  Opt. Classifier:

$$f^*(\mathbf{x}) = \operatorname{argmin}_{j \in \tau} \sum_{k \in \tau} L(j, k) P(t = k|\mathbf{x})$$

$$R^* = E \left[ \min_{j \in \tau} \sum_{k \in \tau} L(j, k) P(t = k|\mathbf{x}) \right]$$

Probab. Models. Max. Likelihood

**MLE**  $\hat{p}_1 = \operatorname{argmax}_{p_1 \in [0,1]} P(D|p_1)$

Maximize  $\log P(D|p_1) : \frac{d \log P(D|p_1)}{dp_1} = 0$  or minimize  $-\log P(D|p_1)$

Gaussian

$$N(x|\mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

Multivariate  $N(\mathbf{x}|\mu, \Sigma) =$

$$|2\pi\Sigma|^{-1/2} e^{-\frac{1}{2}(\mathbf{x}-\mu)^T \Sigma^{-1}(\mathbf{x}-\mu)}$$

Covariance

$$Cov(\mathbf{x}, \mathbf{y}) = E[\mathbf{x}\mathbf{y}^T] - E[\mathbf{x}]E[\mathbf{y}]$$

$$Cov[\mathbf{A}\mathbf{x}] = \mathbf{A}Cov[\mathbf{x}]\mathbf{A}^T$$

Sample

$$\mathbf{S} = \frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})^T = \frac{1}{n} \mathbf{X}^T \mathbf{X}$$

if  $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i = 0$

Correlation

$$\frac{Cov(x_j, x_k)}{\sqrt{Var(x_j)Var(x_k)}} \in [-1, 1]$$

ML plugin discriminant

$$\hat{y}(\mathbf{x}) = \hat{\mathbf{w}}^T \mathbf{x} - \frac{1}{2}(\|\hat{\mu}_{+1}\|^2 - \|\hat{\mu}_{-1}\|^2)$$

$$+ \log \frac{\hat{\pi}_1}{1 - \hat{\pi}_1} \text{ where } \hat{\pi}_1 = n_1/n \text{ and}$$

$$\hat{\mathbf{w}} = \hat{\mu}_{+1} - \hat{\mu}_{-1}$$

$$y_k^*(\mathbf{x}) = -\frac{1}{2} \|\mathbf{x} - \mu_k\|^2 + \log P(t = k) + C$$

Naive Bayes Classifier

$$P(\mathbf{x}|N, t = k) = \prod_{m=1}^M \left( p_m^{(k)} \right)^{\phi_m(\mathbf{x})}$$

with  $\phi_m(\mathbf{x}) = \sum_{j=1}^N I_{\{x_j = m\}}$

Generalization. Regularization

Training error:  
 $\hat{R}_n(f) = \frac{1}{n} \sum_{i=1}^n I_{\{(f(\mathbf{x}_i) \neq t_i)\}}$   
Generalization error:  
 $R(f) = E^* [I_{\{(f(\mathbf{x}) \neq t)\}}]$   
Tikhonov Regularization term  $\frac{\nu}{2} \|\mathbf{w}\|^2 \rightarrow (\Phi^T \Phi + \nu \mathbf{I}) \mathbf{w} = \Phi^T \mathbf{t}$

**MAP**  $\hat{p}_1 = \operatorname{argmax}_{p_1} p(p_1|D) = \operatorname{argmax}_{p_1} P(D|p_1)p(p_1)$

Cond. Likelihood. Logistic Reg.

Maximize Conditional Likelihood  
 $p(\mathbf{t}|\mathbf{w}) = \prod_{i=1}^n p(t_i|y_i)$

Conditional Likelihood Bridge

Likelihood  $p(\mathbf{t}|\theta) \xrightarrow{-\log} E(t, \theta)$  Loss function

**Logistic regression**  $P(t|y) = \sigma(ty)$

$$E_{log}(\mathbf{w}) = -\log P(\mathbf{t}|\mathbf{w})$$
$$= \sum_{i=1}^n \log \left( 1 + e^{-t_i y_i} \right)$$
$$= \sum_{i=1}^n -\log \sigma(t_i y_i)$$

Generative Modeling

$$p(\mathbf{x}, t|\theta) = p(\mathbf{x}|t, \theta)P(t|\theta)$$

joint max. likeli.  $\max_{\theta} \prod_{i=1}^n p(\mathbf{x}_i, t_i|\theta)$

Discriminative Modeling

$$P(t|\mathbf{x}, \theta) \rightarrow \max_{\theta} \prod_{i=1}^n P(t_i|\mathbf{x}_i, \theta)$$

conditional maximum likelihood

Multiway logistic Regression

P(t = k|x) = \frac{e^{y\_k^\*(x)}}{\sum\_{\tilde{k}} e^{y\_{\tilde{k}}^\*(x)}} = \sigma\_k(y^\*(x))

Soft-max mapping \sigma\_k(\nu) = e^{\nu\_k - lsexp(\nu)} with lsexp(\nu) = \log \sum\_{\tilde{k}} e^{\nu\_{\tilde{k}}} \nabla\_v lsexp(v) = \sigma(v)

Support Vector Machines

Kernel function: K(x, x') = \phi(x)^T \phi(x')

Gaussian: e^{-\frac{\tau}{2} ||x - x'||^2}, \tau > 0

Polynomial: K(x, x') = (x^T x')^r

Max margin perceptron

\max\_{w,b} \left\{ \gamma\_D(w,b) = \min\_{i=1..n} \frac{t\_i(w^T \phi(x\_i) + b)}{||w||} \right\}

! D linearly separable !

Hard margin (convex optimization problem): \min\_{w,b} \frac{1}{2} ||w||^2 subj. to t\_i(w^T \phi(x\_i) + b) \ge 1, i = 1..n

Soft margin SVM

\min\_{w,b,\xi} \frac{1}{2} ||w||^2 + C \sum\_{i=1}^n \xi\_i

subj. to t\_i(w^T \phi(x\_i) + b) \ge 1 - \xi\_i, \xi\_i \ge 0

\equiv \min\_{w,b} \frac{1}{2C} ||w||^2 + \underbrace{\sum\_{i=1}^n [1 - t\_i y\_i]\_+}\_{E\_{svm}(w,b)}

Repr. Thm w\_\* = \sum\_{i=1}^n \alpha\_\* i \phi(x\_i)

\rightarrow y(x) = \sum\_{i=1}^n \alpha\_i K(x, x\_i) + b

**Solution** Primal/Dual: p\_\* = \min\_w b, \max\_{0 \le \alpha\_i \le C} L(w, b, \alpha)

d\_\* \max-min (reversed). Weak duality d\_\* \le p\_\* (strong duality iff =)

L(w, b, \alpha) = \frac{1}{2} ||w||^2 + \sum\_{i=1}^n \alpha\_i (1 - t\_i y\_i)

Maximize Criterion for dual:

\phi\_D(\alpha) = \sum\_i \alpha\_i - \frac{1}{2} \sum\_{i,j} \alpha\_i t\_i K\_{ij} t\_j \alpha\_j

subj. to \alpha\_i \in [0, C], \sum\_i \alpha\_i t\_i = 0

Discriminant y^\*(x) = w\_\*^T \phi(x) = \sum\_{i=1}^n \alpha\_{\*,i} t\_i K(x, x\_i) + b\_\*

b = \frac{1}{|S|} \sum\_{i \in S} (t\_i - \hat{y}\_i), S = essential support vectors

Support vectors

\left\{ \begin{array}{ll} \alpha\_i = 0 & 1 - t\_i y\_i \le 0 \text{ not} \\ \alpha\_i \in (0, C) & 1 - t\_i y\_i = 0 \text{ essential} \\ \alpha\_i = C & 1 - t\_i y\_i \ge 0 \text{ bound} \end{array} \right.

Model Selection and Evaluation

Bias-Variance Decomposition

E[(\hat{y}(x|D) - E[t|x])^2|x] =

\underbrace{(E[\hat{y}(x|D)|x] - E[t|x])^2}\_{Bias^2} + \underbrace{Var(\hat{y}(x|D)|x)}\_{Variance}

Ens. meth. \hat{y}\_{ens}(x) = \frac{1}{L} \sum\_{i=1}^L \hat{y}\_i(x)

CV \hat{R}\_{CV}^{(M)}(D) = \frac{1}{n} \sum\_{i=1}^n L(\hat{y}\_{\nu}^{-m(i)}, t\_i)

Dimensionality Reduction

PCA z = U^T x with U^T U = I\_{M \times M}

u\_\* = \operatorname{argmax}\_{u: ||u||=1} u^T Cov(x) u

PC directions = eigendirections of Cov(x)

**Goals:** maximize Cov(z) / Minimize E[||\hat{x} - x||^2] / decorrelate components of z

! PCA doesn't depends on labels t !

Fischer

J(u) = \frac{(m\_1 - m\_0)^2}{s\_0^2 + s\_1^2} = \frac{u^T S\_B u}{u^T S\_W u}

with \mu\_k = \mu^T \mu\_k. Maximize ratio!

S\_B = (\hat{\mu}\_1 - \hat{\mu}\_0)(\hat{\mu}\_1 - \hat{\mu}\_0)^T = dd^T

S\_W = n^{-1} \sum\_{i=1}^n (x\_i - \hat{\mu}\_{t\_i})(x\_i - \hat{\mu}\_{t\_i})^T

\hat{\mu}\_{FLD} = \frac{S\_W^{-1} d}{||S\_W^{-1} d||}, d = \hat{\mu}\_1 - \hat{\mu}\_0

Total covariance S = S\_W + \alpha(1 - \alpha)dd^T with \alpha = \frac{n\_1}{n}

**LDA** Generalization for multiple classes: S = S\_W + S\_B, S\_B = n^{-1} \sum\_k n\_k d\_k d\_k^T, d\_k = \hat{\mu}\_k - \hat{\mu}

\max\_{U \in \mathbb{R}^{D \times M}} tr(U^T S\_B U) s.t. U^T S\_W U = I

Unsupervised Learning

**Kmeans** Iterate until assignment no longer change. Assignment step: ||x\_i - \mu\_{t\_i}|| = \min\_{k=1..K} ||x\_i - \mu\_k||

Update prototypes: \mu\_k = n\_k^{-1} \sum\_{i=1}^n I\_{\{t\_i=k\}} x\_i

\phi(t, \mu) = \sum\_{i=1}^n \sum\_{k=1}^K I\_{\{t\_i=k\}} ||x\_i - \mu\_k||^2

Gaussian Mixture Model

p(x) = \sum\_{k=1}^K p(x|t=k) P(t=k)

= \sum\_{k=1}^K N(x|\mu\_k, \Sigma\_k) P(t=k)

Compute: n\_k = \sum\_{i=1}^n P(t\_i = k|x\_i)

Update: \pi\_k = \frac{n\_k}{n}

\mu\_k = \frac{1}{n\_k} \sum\_{i=1}^n P(t\_i = k|x\_i) x\_i

**Expectation Maximization** E-step: Q\_i(h\_i) \leftarrow P(h\_i|x\_i, \theta)

M-step: maximize surrogate criterion

E(\theta; \{Q\_i\}) = \sum\_{i=1}^n E\_i(\theta; Q\_i)

= \sum\_{i=1}^n E\_{Q\_i}[\log P(x\_i, h\_i|\theta)]

Perfect for missing data ! Hint:

\frac{\partial \log p(x\_i)}{\partial h\_k} = \frac{1}{p(x\_i)} \frac{\partial p(x\_i|h\_k) P(h\_k)}{\partial h\_k}

= \frac{p(x\_i|h\_k) P(h\_k)}{p(x\_i)} \frac{\partial \log p(x\_i|h\_k)}{\partial h\_k}

= P(h\_k|x\_i) \frac{\partial \log p(x\_i|h\_k)}{\partial h\_k}

Beautiful Maths

**Cauchy-Schwarz** |a^T b| \le ||a|| ||b||

**Logistic function** \sigma(v) = \frac{1}{1+e^{-v}}

\sigma'(v) = \sigma(v) \sigma(-v) = \sigma(v) (1 - \sigma(v))

**tanh** g(a) = \tanh(a) = \frac{e^a - e^{-a}}{e^a + e^{-a}}

g(a)' = 1 - g(a)^2

**Trace** tr(A) = \sum\_{j=1}^d a\_{jj} = \sum\_{j=1}^d \lambda\_j

x^T A x = tr(x^T A x) = tr(A x x^T)

**Eigs** A v = \lambda v, A = U \Lambda U^{-1}

|A| = \prod\_{j=1}^d \lambda\_j

**Positive semi-definite matrix** A \in \mathbb{R}^{d \times d} symmetric: v^T A v \ge 0 \forall v \in \mathbb{R}^d, v \ne 0

**Beta**(\alpha, \beta)

p(p\_1|\alpha, \beta) = \frac{1}{B(\alpha, \beta)} (p\_1)^{\alpha-1} (1-p\_1)^{\beta-1}

With B(\alpha, \beta) = \frac{\Gamma(\alpha) \Gamma(\beta)}{\Gamma(\alpha+\beta)}. Mode \frac{\alpha-1}{\alpha+\beta-2}

**Hinge function** [x]\_+ = max(x, 0)

**Cross-entropy, divergence**

D(q||p) = \sum\_{l=1}^L q\_l \log(\frac{q\_l}{p\_l}) \ge 0