

Probabilities

Expectation

$$\begin{aligned}\mathbb{E}[X] &= \int_{\Omega} x f(x) dx = \int_{\omega} x P[X=x] dx \\ \mathbb{E}_{Y|X}[Y] &= \mathbb{E}_Y[Y|X] \\ \mathbb{E}_{X,Y}[f(X,Y)] &= \mathbb{E}_X \mathbb{E}_{Y|X}[f(X,Y)|X] \\ \mathbb{E}_{Y|X}[f(X,Y)|X] &= \int_{\mathbb{R}} f(X,y) p_{Y|X}(y) dy\end{aligned}$$

Variance & Covariance

$$\begin{aligned}\text{Var}(X) &= \mathbb{E}[(X - \mathbb{E}[X])^2] = \mathbb{E}[X^2] - \mathbb{E}[X]^2 \\ \text{Var}[X + Y] &= \text{Var}[X] + \text{Var}[Y] \quad XY \text{ iid} \\ \text{Var}[\alpha X] &= \alpha^2 \text{Var}[X] \\ \text{Cov}(X, Y) &= \mathbb{E}[(X - \mathbb{E}[X])(Y - \mathbb{E}[Y])]\end{aligned}$$

Conditional Probabilities

$$P[X|Y] = \frac{P[X,Y]}{P[Y]}, \quad P[\bar{X}|Y] = 1 - P[X|Y]$$

Distributions

$$\begin{aligned}\mathcal{N}(x|\mu, \sigma^2) &= 1/(\sqrt{2\pi\sigma^2}) \exp^{-(x-\mu)^2/(2\sigma^2)} \\ \mathcal{N}(x|\mu, \Sigma) &= \frac{1}{(2\pi)^{2D/2} |\Sigma|^{1/2}} e^{-\frac{1}{2}(x-\mu)^T \Sigma^{-1} (x-\mu)} \\ \text{Exp}(x|\lambda) &= \lambda e^{-\lambda x} \quad \text{Ber}(x|\theta) = \theta^x (1-\theta)^{(1-x)} \\ \text{Sigmoid: } \sigma(x) &= 1/(1 + \exp(-x))\end{aligned}$$

Chebyshev & Consistency

$$\begin{aligned}P(|X - \mathbb{E}[X]| \geq \epsilon) &\leq \frac{\text{Var}(X)}{\epsilon^2} \\ \lim n \rightarrow \infty P(|\hat{\mu} - \mu| > \epsilon) &= 0\end{aligned}$$

Cramer Rao lower bound

$$\begin{aligned}\text{Var}[\hat{\theta}] &\geq \mathcal{I}_n(\theta) \\ \mathcal{I}_n(\theta) &= -\mathbb{E}\left[\frac{\partial^2 \log[\mathcal{L}_n|\theta]}{\partial \theta^2}\right] \quad \hat{\theta} \text{ unbiased} \\ \text{Efficiency of } \hat{\theta}: e(\theta_n) &= \frac{1}{\text{Var}[\hat{\theta}_n] \mathcal{I}_n(\theta)}\end{aligned}$$

$$\begin{aligned}e(\theta_n) &= 1 \text{ (efficient)} \\ \lim_{n \rightarrow \infty} e(\theta_n) &= 1 \text{ (asympt. efficient)}\end{aligned}$$

Matrix Derivations

$$\begin{aligned}\frac{\partial \mathbf{a}^T \mathbf{x}}{\partial \mathbf{x}} &= \mathbf{a} \quad \frac{\partial \mathbf{a}^T \mathbf{x} \mathbf{b}}{\partial \mathbf{x}} = \mathbf{a} \mathbf{b}^T \quad \frac{\partial \mathbf{a}^T \mathbf{x}^T \mathbf{b}}{\partial \mathbf{x}} = \mathbf{b} \mathbf{a}^T \\ \frac{\partial \mathbf{a}^T \mathbf{X} \mathbf{a}}{\partial \mathbf{x}} &= \mathbf{a}^T (\mathbf{X} + \mathbf{X}^T) \\ \frac{\partial \mathbf{f}(\mathbf{x})^T \mathbf{g}(\mathbf{x})}{\partial \mathbf{x}} &= \mathbf{f}(\mathbf{x})^T \frac{\partial \mathbf{g}(\mathbf{x})}{\partial \mathbf{x}} + \mathbf{g}(\mathbf{x})^T \frac{\partial \mathbf{f}(\mathbf{x})}{\partial \mathbf{x}}\end{aligned}$$

$\mathbf{X}^T \mathbf{X}$: only invertible if none of the Eigenvalue is 0. Inversion instable if ratio from \mathbf{X} 's smallest EV to the largest is big.

Optimization

Gradient Descent

$$\begin{aligned}\theta^{\text{new}} &\leftarrow \theta^{\text{old}} - \eta \nabla_{\theta} \mathcal{L} \\ \text{Convergence isn't guaranteed.} \\ \text{Less zigzag by adding momentum:} \\ \theta^{(l+1)} &\leftarrow \theta^{(l)} - \eta \nabla_{\theta} \mathcal{L} + \mu(\theta^{(l)} - \theta^{(l-1)})\end{aligned}$$

Newton's Method

$$\begin{aligned}\text{Use 2nd order derivation. (Hessian)} \\ \theta^{\text{new}} &\leftarrow \theta^{\text{old}} - \eta (\nabla_{\theta} \mathcal{L} / \nabla_{\theta}^2 \mathcal{L}) \\ H = \nabla_{\theta}^2 \mathcal{L} &\text{ has to be p.d (convex func).}\end{aligned}$$

Risks and Losses

Expected Risk

$$\begin{aligned}\text{Conditional Expected Risk} \\ R(f, X) &= \int_{\mathbb{R}} \mathcal{L}(Y, f(X)) P(Y|X) dY \\ \text{Total Expected Risk } R(f) &= \mathbb{E}_X[R(f, X)] = \int_{\mathcal{X}} R(f, X) P(X) dX = \int_{\mathcal{X}} \int_{\mathbb{R}} \mathcal{L}(Y, f(X)) P(X, Y) dX dY\end{aligned}$$

Empirical Risk

$$\begin{aligned}Z^{\text{train}} &= (X_1, Y_1), \dots, (X_n, Y_n) \\ Z^{\text{test}} &= (X_{n+1}, Y_{n+1}), \dots, (X_{n+m}, Y_{n+m}) \\ \text{Empirical Risk Minimizer } \hat{f} &\text{ s.t.} \\ \hat{f} &\in \arg \min_{f \in \mathcal{C}} \hat{R}(\hat{f}, Z^{\text{train}})\end{aligned}$$

Training error:

$$\begin{aligned}\hat{R}(\hat{f}, Z^{\text{train}}) &= \frac{1}{n} \sum_{i=1}^n Q(Y_i, \hat{f}(X_i)) \\ \text{Test error:} \\ \hat{R}(\hat{f}, Z^{\text{test}}) &= \frac{1}{m} \sum_{i=n+1}^{n+m} Q(Y_i, \hat{f}(X_i)) \\ \hat{R}(\hat{f}, Z^{\text{test}}) &\neq \mathbb{E}_X[R(f, X)]\end{aligned}$$

Linear Regression

$$\begin{aligned}\text{Data: } Z &= (x_i, y_i) \in \mathbb{R}^3 \times \mathbb{R} : 1 \leq i \leq n \\ X &\text{ are iids and } Y \text{ depends on } X. \\ \text{Model: } \mathbf{Y} &= \beta_0 + \sum_{j=1}^d \mathbf{X}_j \beta_j \quad \mathbf{Y} \subset \mathbb{R} \\ \text{Introduce } X_0 &= 1 \text{ and rewrite} \\ \mathbf{Y} = \mathbf{X}^T \boldsymbol{\beta} \quad \mathbf{X} &\in \mathbb{R}^{(d+1) \times n}, \boldsymbol{\beta} \in \mathbb{R}^{d+1} \\ \text{additive Gaussian noise } \epsilon &\sim \mathcal{N}(0, \sigma^2)\end{aligned}$$

$$\begin{aligned}\hat{y} &= \mathbf{X} \hat{\boldsymbol{\beta}} + \epsilon \\ \hat{\boldsymbol{\beta}} &\sim \mathcal{N}(\boldsymbol{\beta}, (\mathbf{X}^T \mathbf{X})^{-1} \sigma^2) \text{ and} \\ p(Y|\mathbf{X}, \boldsymbol{\beta}, \sigma) &\sim \mathcal{N}(Y|\mathbf{X}^T \boldsymbol{\beta}, \sigma^2) \\ \text{A Regression has Optimum:} \\ f^*(x) &= \mathbb{E}_Y[Y|X=x]\end{aligned}$$

Linear Regression

$$\begin{aligned}\text{Setting: Minimize RSS.} \\ \mathcal{L} &= \text{RSS}(\boldsymbol{\beta}) = \sum_{i=1}^n (y_i - x_i^T \boldsymbol{\beta})^2 = \\ &= (\mathbf{y} - \mathbf{X} \boldsymbol{\beta})^T (\mathbf{y} - \mathbf{X} \boldsymbol{\beta}) \\ X &\in \mathbb{R}^{n \times (d+1)}, y \in \mathbb{R}^n, \boldsymbol{\beta} \in \mathbb{R}^{d+1} \\ \text{Solution: differentiate w.r.t } \boldsymbol{\beta} \\ \hat{\boldsymbol{\beta}} &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \\ \text{Is an orth. projection with lowest variance of all unbiased estimates.} \\ \text{Prediction: } \hat{y} &= \mathbf{X} \hat{\boldsymbol{\beta}} = \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}\end{aligned}$$

Ridge Regression (L2 penalty)

$$\begin{aligned}\text{Setting: Penalize the } \boldsymbol{\beta}\text{s} \\ \mathcal{L} &= \sum_{i=1}^n (y_i - x_i^T \boldsymbol{\beta})^2 + \lambda \sum_{j=1}^d \beta_j^2 = \\ &= (\mathbf{y} - \mathbf{X} \boldsymbol{\beta})^T (\mathbf{y} - \mathbf{X} \boldsymbol{\beta}) + \lambda \boldsymbol{\beta}^T \boldsymbol{\beta} \\ \text{Solution: differentiate w.r.t } \boldsymbol{\beta} \\ \hat{\boldsymbol{\beta}}^{\text{ridge}} &= (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y}\end{aligned}$$

Lasso (L1 penalty)

$$\begin{aligned}\text{Setting: seek for a sparse solution} \\ \mathcal{L} &= \sum_{i=1}^n (y_i - x_i^T \boldsymbol{\beta})^2 + \lambda \sum_{j=1}^d |\beta_j|\end{aligned}$$

$$\begin{aligned}&= (\mathbf{y} - \mathbf{X} \boldsymbol{\beta})^T (\mathbf{y} - \mathbf{X} \boldsymbol{\beta}) + \lambda \|\boldsymbol{\beta}\|_1 \\ \text{Lasso has no closed form.} \\ \text{Bayesian Linear Regression} \\ \text{Setting: Define a prior over the } \boldsymbol{\beta}\text{s.} \\ \text{e.g. Ridge:} \\ \text{Assume } \boldsymbol{\beta}\text{s distributed with mean 0} \\ p(\boldsymbol{\beta}|\boldsymbol{\Lambda}) &= \mathcal{N}(\boldsymbol{\beta}|\mathbf{0}, \boldsymbol{\Lambda}^{-1}) \propto \exp(-\frac{1}{2} \boldsymbol{\beta}^T \boldsymbol{\Lambda} \boldsymbol{\beta}) \\ \text{e.g. Linear Regression:} \\ \text{equivalent to ridge with } \boldsymbol{\Lambda} &= \lambda \mathbf{I}, \sigma = 1\end{aligned}$$

Posterior

$$\begin{aligned}\text{given observed } \mathbf{X}, \mathbf{y}, \text{ use Baye's theorem to find the posterior} \\ p(\boldsymbol{\beta}|\mathbf{X}, \mathbf{y}, \boldsymbol{\Lambda}, \sigma) &= \mathcal{N}(\boldsymbol{\mu}_{\boldsymbol{\beta}}, \boldsymbol{\Sigma}_{\boldsymbol{\beta}}) \\ \boldsymbol{\mu}_{\boldsymbol{\beta}} &= \sigma^2 (\mathbf{X}^T \mathbf{X} + \sigma^2 \boldsymbol{\Lambda})^{-1} (\mathbf{X})^T \mathbf{y} \\ \boldsymbol{\Sigma}_{\boldsymbol{\beta}} &= \sigma^2 (\mathbf{X}^T \mathbf{X} + \sigma^2 \boldsymbol{\Lambda})^{-1}\end{aligned}$$

Nonlinear Regression

$$\begin{aligned}\text{Idea: Feature space transformation} \\ \text{Model: } \mathbf{Y} = f(\mathbf{X}) &= \sum_{m=1}^M \beta_m h_m(\mathbf{X}) \\ \text{Transformation } h_m(\mathbf{X}): \mathbb{R}^d &\rightarrow \mathbb{R}\end{aligned}$$

Cubic Spline

$$\begin{aligned}\text{e.g. for } d=1 \text{ with knots at } \xi_1 \text{ and } \xi_2 \\ h_1(X) &= 1 \quad h_3(X) = X^2 \quad h_5(X) = (X - \xi_1)_+^3 \\ h_2(X) &= X \quad h_4(X) = X^3 \quad h_6(X) = (X - \xi_2)_+^3\end{aligned}$$

Wavelets

Functions that measure local properties of the underlying data. Keep the most important ones and get rid of the noise.

Gaussian Process Regression

$$\begin{aligned}\text{joint Gaussian over all outputs} \\ \mathbf{y} = \mathbf{X} \boldsymbol{\beta} + \epsilon \quad \epsilon &\sim \mathcal{N}(\epsilon|\mathbf{0}, \sigma^2 \mathbf{I}_n) \\ \text{We can rewrite the distribution} \\ P\left(\begin{bmatrix} \mathbf{y} \\ \mathbf{y}_* \end{bmatrix}\right) &= \mathcal{N}(\mathbf{y}|\mathbf{0}, \begin{bmatrix} \mathbf{C}_n & \mathbf{k} \\ \mathbf{k}^T & c \end{bmatrix})\end{aligned}$$

Such that for prediction:

$$\begin{aligned}p(y_*|\mathbf{x}_*, \mathbf{X}, \mathbf{y}) &= \mathcal{N}(y_*|\boldsymbol{\mu}_*, \sigma_*^2) \\ \boldsymbol{\mu}_{y_*} &= \mathbf{k}^T \mathbf{C}_n^{-1} \mathbf{y} \quad \mathbf{C}_n = \mathbf{K} + \sigma^2 \mathbf{I} \\ \sigma_*^2 &= c - \mathbf{k}^T \mathbf{C}_n^{-1} \mathbf{k} \quad c = k(x_*, x_*) + \sigma^2 \\ \mathbf{k} &= k(x_*, \mathbf{X})\end{aligned}$$

k is the kernel function. lengthscale in kernel: how far can we reliably extrapolate

Bias-Variance tradeoff

$$\begin{aligned}\text{Bias}(\hat{f}) &= \mathbb{E}[\hat{f}] - f^* \\ \text{Var}(\hat{f}) &= \mathbb{E}[(\hat{f} - \mathbb{E}[\hat{f}])^2] \\ |\mathcal{Z}| \downarrow \quad |\mathcal{F}| \uparrow &\Rightarrow \text{Var} \uparrow \quad \text{Bias} \downarrow \\ |\mathcal{Z}| \uparrow \quad |\mathcal{F}| \downarrow &\Rightarrow \text{Var} \downarrow \quad \text{Bias} \uparrow\end{aligned}$$

Squared Error Decomposition

$$\begin{aligned}\mathbb{E}_D \mathbb{E}_{X,Y}[(\hat{f}(X) - Y)^2] &= \\ \mathbb{E}_{X,Y}[(\mathbb{E}_Y[Y|X] - Y)^2] &\text{ (noise)}\end{aligned}$$

$$\begin{aligned}&+ \mathbb{E}_X \mathbb{E}_D[(\hat{f}_D(X) - \mathbb{E}_D[\hat{f}(X)])^2] \text{ (var.)} \\ &+ \mathbb{E}_X[(\mathbb{E}_D[\hat{f}_D(X)] - \mathbb{E}_Y[Y|X])^2] \text{ (bias}^2\text{)} \\ &\text{(can be derivated by vanishing of the crossproducts)}\end{aligned}$$

Parametric Density Estimation

$$\begin{aligned}\text{Find the most likely parameter of a distribution.} \\ \text{Maximum Likelihood} \\ \text{Likelihood: } P(\mathcal{X}|\theta) &= \prod_{i \leq n} p(x_i|\theta) \\ \text{Find: } \hat{\theta} &\in \arg \max_{\theta} P(\mathcal{X}|\theta) \\ \text{Procedure: solve } \nabla_{\theta} \log P(\mathcal{X}|\theta) &= 0 \\ \text{Efficient \& easy to calculate.} \\ \text{Consistent. Converge to best model } \theta_0 &\text{ Warning: Overfitting!}\end{aligned}$$

Maximum A Posteriori

$$\begin{aligned}\text{Assume Knowledge of a prior } P(\theta) \\ \text{Find: } \hat{\theta} &\in \arg \max_{\theta} P(\theta|\mathcal{X}) = \\ &= \arg \max_{\theta} P(\mathcal{X}|\theta) P(\theta) \\ \text{Solve } \nabla_{\theta} \log P(\mathcal{X}|\theta) P(\theta) &= 0\end{aligned}$$

0.1 Bayesian Learning

$$\begin{aligned}\text{Prior Knowledge of } p(\theta) \\ \text{Find Posterior Density: } p(\theta|\mathcal{X}) \\ \text{Can be done using Baye's Rules} \\ \text{We can use this Recursively:}\end{aligned}$$

$$\begin{aligned}\mathcal{X}^n &= \{x_1, \dots, x_n\} \\ p(\theta|\mathcal{X}^n) &= \frac{p(x_n|\theta) p(\theta|\mathcal{X}^{n-1})}{\int p(x_n|\theta) p(\theta|\mathcal{X}^{n-1}) d\theta} \text{ with}\end{aligned}$$

$$\begin{aligned}p(\theta|\mathcal{X}^0) p(\theta) \\ \text{Difficult \& needs prior knowledge.} \\ \text{But better against overfitting.}\end{aligned}$$

Numerical Est. Techniqes

$$\begin{aligned}\text{Setting: Estimate } \hat{f}(x) \in \mathcal{F} \text{ with minimal prediction error.}\end{aligned}$$

K-Fold Cross Validation

$$\begin{aligned}\text{Initialisation (split training set):} \\ \mathcal{Z} = \mathcal{Z}_1 \cup \mathcal{Z}_2 \cup \dots \cup \mathcal{Z}_K, \mathcal{Z}_\mu \cap \mathcal{Z}_\nu &= \emptyset \\ \text{with map } \kappa: \{1, \dots, n\} &\rightarrow \{1, \dots, K\} \\ |\mathcal{Z}_k| &\approx n \frac{K-1}{K} \\ \text{Learning:}\end{aligned}$$

$$\hat{f}^{-\nu}(x) = \arg \min_{f \in \mathcal{F}} \frac{\sum_{i \in \mathcal{Z}_\nu} (y_i - f(x_i))^2}{|\mathcal{Z} - \mathcal{Z}_\nu|}$$

$$\begin{aligned}\text{Validation:} \\ \hat{R}^{cv} &= \frac{1}{n} \sum_{i \leq n} (y_i - \hat{f}^{-\kappa(i)}(x_i))^2 \\ \text{tendance to Underfit} \\ \text{Leave-one-out: } K &= n \text{ (unbiased but Var can be large } \leftarrow \text{corr. datasets)}\end{aligned}$$

Bootstrapping

$$\begin{aligned}\text{Bootstrap samples: } \mathcal{Z}^* &= \{\mathcal{Z}_1^*, \dots, \mathcal{Z}_n^*\} \\ \text{each data point in } \mathcal{Z}_i^* &\text{ was randomly drawn from } \mathcal{Z} \text{ with replacement.} \\ e_0 \text{ Estimator: the error rate for the} &\text{ test data (data that wasn't selected by}\end{aligned}$$

the bootstrap) is assumed to be the error estimate (e.g. for classification):

$$\hat{R}(S(\mathcal{Z})) = \frac{1}{B} \sum_{b=1}^B \sum_{z_i \notin \mathcal{Z}^{*b}} \frac{\mathbb{I}_{c(z_i) \neq y_i}}{|n - \mathcal{Z}^{*b}|}$$

Jackknife

$$\begin{aligned}\text{Estimate of an Estimator } \hat{S}_n \text{'s Bias.} \\ \hat{S}^{JK} &= \hat{S}_n - \text{bias}^{JK} \text{ is JK Estimator.} \\ \text{bias}^{JK} &= (n-1)(\hat{S}_n - \hat{S}_m) \\ \hat{S}_n &= \frac{1}{n} \sum_{i=1}^n \hat{S}_{n-1}^{(-i)} \text{ avg. LOO Estimator} \\ \text{Debiased est. can have big variance!} \\ \text{Bootstrap} &\quad \text{Debiased} \\ \bar{S} &= 2\hat{S} - \frac{1}{B} \sum_b \hat{S}^*(b)\end{aligned}$$

Classification

$$\begin{aligned}\text{group points in classes } 1, \dots, k, \mathcal{D}, \mathcal{O} \\ \mathcal{D}: \text{doubt class, } \mathcal{O}: \text{outliers.} \\ \text{Data: } \mathcal{Z} = \{z_i = (x_i, y_i) : 1 \leq i \leq n\} &\text{ Assume we know } p_y(x) = P[X=x|Y=y] \\ \text{Found: classifier } \hat{c}: \mathcal{X} \rightarrow \mathcal{Y} &= \{1, \dots, \mathcal{D}\} \\ \text{Error: } \hat{R}(\hat{c}|\mathcal{Z}) &= \sum_{(x_i, y_i) \in \mathcal{Z}} \mathbb{I}_{\{\hat{c}(x_i) \neq y_i\}} \\ \text{Expected Error:} \\ \mathcal{R}(\hat{c}) &= \sum_{y \leq k} P[y] \mathbb{E}_{x|y}[\mathbb{I}_{\{\hat{c}(x_i) \neq y_i\}} | Y=y] \\ &\text{(add term from } \mathcal{D})\end{aligned}$$

Loss Functions

$$\begin{aligned}0\text{-1 Loss: } L^{0-1}(y, z) &= \begin{cases} 0 & \text{if } (z = y) \\ 1 & \text{if } (z \neq y) \end{cases}\end{aligned}$$

$$\begin{aligned}\text{Exponential Loss:} \\ L^{\text{exp}}(y, z) &= \exp(-(2y-1)(2z-1)) \\ \text{Logistic Loss:} \\ L^{\text{log}}(y, z) &= \ln(1 + \exp((2y-1)(2z-1))) \\ \text{Hinge Loss:} \\ \text{Favors sparsity. Used in SVM} \\ L^{\text{hinge}}(y, z) &= \max\{0, 1 - (2y-1)(2z-1)\}\end{aligned}$$

Bayes Optimal Classifier

$$\begin{aligned}\text{Minimizes total risk for 0-1 Loss} \\ \hat{c}(x) = \begin{cases} y & \text{if } p(y|x) = \max_{z \leq k} p(z|x) > 1 - \alpha \\ \mathcal{D} & \text{if } p(y|x) < 1 - \alpha \end{cases} \end{aligned}$$

Generalize to other loss functions

Discriminant Functions

$$\begin{aligned}\text{Functions } g_k(x) \quad 1 \leq k \leq K \\ \text{Decide: } g_y(x) > g_z(x) \forall z \neq y \Rightarrow \text{chose } y \\ \text{Const factor doesn't change decision.} \\ g_k(x) &= P[y|x] \propto P[x|y] P[y] \Rightarrow \\ g_k(x) &= \ln P[x|y] + \ln P[y] = \ln P[x|y] + \pi_y \\ \text{implements an opt. Baye classifier.}\end{aligned}$$

Decision Surface of Discriminant

$$\begin{aligned}\text{Solve: } g_{k_1}(x) - g_{k_2}(x) = 0 \text{ Special case with Gaussian classes:} \\ \text{if } \Sigma_y = \Sigma \Rightarrow \text{linear decision surface} \\ g_k(x) &= w^T (x - x_0) \quad w = \Sigma^{-1} (\mu_1 - \mu_2) \\ x_0 &= \frac{1}{2} (\mu_1 + \mu_2) - \frac{\sigma^2 (\mu_1 - \mu_2)}{(\mu_1 - \mu_2)^T \Sigma^{-1} (\mu_1 - \mu_2)} \log \frac{\pi_1}{\pi_2}\end{aligned}$$

Linear Classifier
 optimal for Gaussian with equal cov.
 Stat. simplicity & comput. efficiency.

$g(x) = a^T \tilde{x}$ $a = (w_0, w)^T, \tilde{x} = (1, x)^T$
 $a^T \tilde{x}_i > 0 \Rightarrow y_i = 1$ $a^T \tilde{x}_i < 0 \Rightarrow y_i = 2$
 Normalization: $\tilde{x}_i \rightarrow -\tilde{x}_i$ if $y_i = 2$
 Find $a: a^T \tilde{x} > 0$ (linearly separable)
 Learning w. Gradient Descent:
 $a(k+1) = a(k) - \eta(k) \nabla J(a(k))$
 $J(\cdot)$: cost function $\eta(\cdot)$: learning rate
 Newton's rule (opt. grad descent):
 $a(k+1) = a(k) - H^{-1} \nabla J$ $H = \frac{\partial^2 J}{\partial a_i \partial a_j}$

Perceptron Criterion
 $J_P(a) = \sum_{\tilde{x} \in \tilde{\mathcal{X}}} (-a^T \tilde{x})$
 $\tilde{\mathcal{X}}$ set of misclassified samples.
 $\Rightarrow a(k+1) = a(k) + \eta(k) \sum_{\tilde{x} \in \tilde{\mathcal{X}}} \tilde{x}$ Converges if data separable.

WINNOW Algorithm
 Performs better when many dimensions are irrelevant. Search for 2 weight vectors a^+, a^- (for each class). If a point is misclassified: $a_i^+ \leftarrow -\alpha^{\tilde{x}_i} a_i^+, a_i^- \leftarrow -\alpha^{-\tilde{x}_i} a_i^-$ (class 1 err.)
 $a_i^+ \leftarrow -\alpha^{-\tilde{x}_i} a_i^+, a_i^- \leftarrow -\alpha^{\tilde{x}_i} a_i^-$ (class 2 err.)
 Exponential update.

Fisher's Linear Discriminant Analysis
 Maximize distance of the means of the projected classes to find projection plane separating them best.
 proj mean: $\tilde{\mu}_\alpha = \frac{1}{n_\alpha} \sum_{x \in \mathcal{X}_\alpha} w^T x = w^T \mu_\alpha$

Dist of proj means: $|w^T (\mu_1 - \mu_2)|$ Classifies proj. cov: $\tilde{\Sigma}_1 + \tilde{\Sigma}_2 = w^T (\Sigma_1 + \Sigma_2) w$
 Fishers Criterion:
 $J(w) = \frac{\|\tilde{\mu}_1 - \tilde{\mu}_2\|^2}{\tilde{\Sigma}_1 + \tilde{\Sigma}_2} = \frac{w^T (\mu_1 - \mu_2)(\mu_1 - \mu_2)^T w}{w^T (\Sigma_1 + \Sigma_2) w}$
 Fishers Crit for Multiple Classes:

$J(W) = \frac{|W^T S_B W|}{W^T S_W W}$
 $S_B = \sum_{i=1}^k n_k (\mu_k - \mu)(\mu_k - \mu)^T$
 $S_W = \sum_{i=1}^k \sum_{x \in \mathcal{D}_i} (x - \mu_i)(x - \mu_i)^T$

Linear Discriminant for Multiclass
 Reformulate as $(k-1)$ "class α - not class α " dichotomie. But some area are ambiguous

Support Vector Machine (SVM)
 Generalize Perceptron with margin and kernel. Find plane that maximizes margin m s.t.

$$z_i g(\mathbf{y}) = z_i (\mathbf{w}^T \mathbf{y} + w_0) \geq m \quad \forall \mathbf{y}_i \in \mathcal{Y}$$

$z_i \in \{-1, +1\}$ $\mathbf{y}_i = \phi(\mathbf{x}_i)$
 Vectors \mathbf{y}_i are the support vectors
 Functional Margin Problem:
 minimizes $\|\mathbf{w}\|$ for $m=1: L(\mathbf{w}, w_0, \alpha) = \frac{1}{2} \mathbf{w}^T \mathbf{w} - \sum_{i=1}^n \alpha_i [z_i (\mathbf{w}^T \mathbf{y}_i + w_0) - 1]$
 where α s are Lagrange multipliers.
 $\frac{\partial L}{\partial w} = 0$ and $\frac{\partial L}{\partial w_0} = 0$ give us constraints
 $\mathbf{w} = \sum_{i=1}^n \alpha_i z_i \mathbf{y}_i$ $0 = \sum_{i=1}^n \alpha_i z_i$
 Replacing these in $L(\mathbf{w}, w_0, \alpha)$ we get
 $\tilde{L}(\alpha) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n \alpha_i \alpha_j z_i z_j \mathbf{y}_i^T \mathbf{y}_j$
 with $\alpha_i \geq 0$ and $\sum_{i=1}^n \alpha_i z_i = 0$
 This is the dual representation. The optimal hyperplane is given by $\mathbf{w}^* = \sum_{i=1}^n \alpha_i^* z_i \mathbf{y}_i$
 $w_0^* = -\frac{1}{2} (\min_{z_i=1} \mathbf{w}^{*T} \mathbf{y}_i + \max_{z_i=-1} \mathbf{w}^{*T} \mathbf{y}_i)$
 where α maximize the dual problem.
 Only Support Vectors ($\alpha_i \neq 0$) contribute to the evaluation.
 Optimal Margin: $\mathbf{w}^T \mathbf{w} = \sum_{i \in S_V} \alpha_i^*$
 Discrim.: $g^*(\mathbf{x}) = \sum_{i \in S_V} z_i \alpha_i \mathbf{y}_i^T \mathbf{y}_i + w_0^*$
 class = sign($\mathbf{y}^T \mathbf{w}^* + w_0^*$)

Soft Margin SVM
 Introduce slack to relax constraints
 $z_i (\mathbf{w}^T \mathbf{y}_i + w_0) \geq m(1 - \xi_i)$
 $L(\mathbf{w}, w_0, \xi, \alpha, \beta) = \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_{i=1}^n \xi_i - \sum_{i=1}^n \alpha_i [z_i (\mathbf{w}^T \mathbf{y}_i + w_0) - 1 + \xi_i] - \sum_{i=1}^n \beta_i \xi_i$
 C controls margin maximization vs. constraint violation
 Dual Problem same than usual SVM but with supplementary constraint: $\alpha_i \leq C$

Non-Linear SVM
 use kernel in discriminant funct:
 $g(\mathbf{x}) = \sum_{i=1}^n \alpha_i z_i K(\mathbf{x}_i, \mathbf{x})$
 E.g solve the XOR Problem with:
 $K(x, y) = (1 + x_1 y_1 + x_2 y_2)^2$

Multiclass SVM
 \forall class $z \in \{1, 2, \dots, M\}$ we introduce \mathbf{w}_z and define the margin m s.t.:
 $(\mathbf{w}_z^T \mathbf{y}_i + w_{z,0}) - \max_{z \neq Z_i} (\mathbf{w}_z^T \mathbf{y}_i + w_{z,0}) \geq 0 \quad \forall \mathbf{y}_i \in \mathcal{Y}$

Structured SVM
 Each sample \mathbf{y} is assigned to a structured output label z
 Output Space Representation:
 joint feature map: $\psi(z, \mathbf{y})$
 scoring function: $f_w(z, \mathbf{y}) = \mathbf{w}^T \psi(z, \mathbf{y})$
 Classify: $\hat{z} = h(\mathbf{y}) \arg \max_{z \in \mathcal{K}} f_w(z, \mathbf{y})$

Kernels
 Similarity based reasoning
 Gram Matrix $K = (K(\mathbf{x}_i, \mathbf{x}_j)) \quad 1 \leq i, j \leq n$
 $K(\mathbf{x}, \mathbf{x}') = \phi(\mathbf{x})^T \phi(\mathbf{x}')$ $K(\mathbf{x}, \mathbf{x}') = K(\mathbf{x}', \mathbf{x})$
 $K(\mathbf{x}, \mathbf{x}')$ pos.semi-def. (all EV ≥ 0)
 If K_1 & K_2 are kernels K is too:
 $K(\mathbf{x}, \mathbf{x}') = K_1(\mathbf{x}, \mathbf{x}') K_2(\mathbf{x}, \mathbf{x}')$
 $K(\mathbf{x}, \mathbf{x}') = \alpha K_1(\mathbf{x}, \mathbf{x}') + \beta K_2(\mathbf{x}, \mathbf{x}')$
 $K(\mathbf{x}, \mathbf{x}') = K_1(h(\mathbf{x}), h(\mathbf{x}')) \quad h: \mathcal{X} \rightarrow \mathcal{X}'$
 $K(\mathbf{x}, \mathbf{x}') = h(K_1(\mathbf{x}, \mathbf{x}')) \quad h$: poly/exp
 Kernel Function Examples:
 $K(\mathbf{x}, \mathbf{x}') = \mathbf{x}^T \mathbf{x}'$ $K(\mathbf{x}, \mathbf{x}') = (\mathbf{x}^T \mathbf{x}' + 1)^p$
 RBF (Gauss): $K(\mathbf{x}, \mathbf{x}') = \exp(-\|\mathbf{x} - \mathbf{x}'\|_2^2 / h^2)$
 Sigmoid: $K(\mathbf{x}, \mathbf{x}') = \tanh(\alpha \mathbf{x}^T \mathbf{x}' + c)$
 not p.s-d eg: $\mathbf{x} = [1, -1], \mathbf{x}' = [-1, 2]$

Ensemble Methods
Combining Regressors
 set of estimators: $\hat{f}_1(x), \dots, \hat{f}_B(x)$ simple average: $\hat{f}(x) = \frac{1}{B} \sum_{i=1}^B \hat{f}_i(x)$
 Bias [$\hat{f}(x)$] = $\frac{1}{B} \sum_{i=1}^B \text{Bias}[f_i(x)]$
 $\text{Var}[\hat{f}(x)] \approx \frac{\sigma}{B}$ if the estimators are uncorrelated.

Combining Classifiers
 Input: classifiers $c_1(x), \dots, c_B(x)$
 Infer $\hat{c}_B(x) = \text{sgn}(\sum_{b=1}^B \alpha_b c_b(x))$
 with weights $\{\alpha_b\}_{b=1}^B$
 Requires diversity of the classifiers.

Bagging
 Train on bootstrapped subsets.
 Sample: $\mathcal{Z} = \{(x_1, y_1), \dots, (x_n, y_n)\}$
 \mathcal{Z}^* : chose i.i.d from \mathcal{Z} w. replacement

Random Forest (Bagging strategy)
 Collection of uncorr. decision trees. Partition data space recursively. Grow the tree sufficiently deep to reduce bias. Prediction with voting.

Boosting
 Combine uncorr. weak learners in sequence. (Weak to avoid overfitting). Coeff. of \hat{c}_{b+1} depend on \hat{c}_b 's results
AdaBoost (minimizes exp. loss)

Init: $\mathcal{X} = \{(x_1, y_1), \dots, (x_n, y_n)\}, w_i^{(1)} = \frac{1}{n}$
 Fit $\hat{c}_b(x)$ to \mathcal{X} weighted by $w^{(b)}$
 $\epsilon_b = \sum_{i=1}^n w_i^{(b)} \mathbb{I}_{\{c_b(x_i) \neq y_i\}} / \sum_{i=1}^n w_i^{(b)}$
 $\alpha_b = \log \frac{1 - \epsilon_b}{\epsilon_b}$
 $w_i^{(b+1)} = w_i^{(b)} \exp(\alpha_i \mathbb{I}_{\{c_b(x_i) \neq y_i\}})$
 return $\hat{c}_B(x) = \text{sgn}(\sum_{b=1}^B \alpha_b c_b(x))$
 best approx. at log-odds ratio.

Neural Networks
Multi Layer Perceptron
 $\{x_j\}_{j=1}^J$ input, $\{y_i\}_{i=1}^I$ output
 $\{z_k^l\}_{k=1}^{K(l)}$ hidden nodes in layer $l \quad 1 \leq l \leq L$
 w_{mk}^l weights from z_k^{l-1} to z_m^l
 $w_i^l k^{L+1}$ weights from z_k^L to output y_i
 $z_k^l = h(a_k^l) = h(\sum_{m=1}^{K(l-1)} w_{km}^l z_m^{l-1})$
 $y_i = \sigma(a_i^{L+1}) = h(\sum_{m=1}^{K(L)} w_{im}^L z_m^L)$
 $\mathcal{L}(\hat{y}(\mathbf{W}, \mathbf{X}), y) = \sum_{n=1}^N \mathcal{L}_n(\hat{y}(\mathbf{W}, \mathbf{X}_n), Y_n)$
 $L = 0$ or $h(a) = a \Rightarrow$ multiple lin. reg.
 Layers \Rightarrow generaliz. & simplicity.
 Model data generating mechanism.
Backpropagation
 Effic. evaluation of loss derivative:
 $\frac{\partial \mathcal{L}_n}{\partial w_{ik}^{l+1}} = \delta_i^{L+1} z_k^l \quad \frac{\partial \mathcal{L}_n}{\partial w_{mk}^l} = \delta_m^l z_k^{l-1}$
 $\delta_i^{L+1} = (\hat{y}_i - y_i) \sigma'(\sum_{m=1}^{K(L)} w_{im}^{L+1} z_m^L)$
 $\delta_m^l = (\sum_{r=1}^{K(l+1)} \delta_r^{l+1} w_{rm}^{l+1}) \cdot h'(\sum_{r=1}^{K(l-1)} w_{mr}^{l-1} z_r^{l-1})$
 $w_{ij}^l \leftarrow w_{ij}^l + \eta \delta_i^l z_j^{(l-1)}$
Regularization
 Avoid overfitting on complex nets.
Early Stopping separate data into train/error/validation sets.
Drop Out Combine thinned nets with removed nodes.
Bayesian priors on w 's
Autoencoder
 Data compression purposes, Output should reproduce input. \Rightarrow PCA
Convolutional Neural Network
 Modelling invariance. Convolutional Layers (filters on a region) & Pooling Layers (aggregate nodes together).
Boltzmann Machine
 Symmetric coupling. Visible and Hidden units. Update through voting of neighbors. Find the weights which generate a defined activity of visible nodes.
Unsupervised Learning
Histograms
 $p_i = \frac{n_i}{N \Delta_i} \quad n \leq N$ in bin i of size Δ_i
 Not scaling to multiple dimensions.
 $K \simeq NP \quad P \simeq p(x) V \Rightarrow p(x) = \frac{K}{N V}$
 K #samples in region of volume V , P probability of falling in it.
Kernel Density Estimator
 Fix V and determine K .

Kernel Window
 $K = \sum_{n=1}^N \phi(\frac{x-x_n}{h}), \phi(u) = \mathbb{I}_{\{\|x\| \leq \frac{1}{2}\}}$
 $p(x) = \frac{1}{N} \sum_{n=1}^N \frac{1}{h^D} \phi(\frac{x-x_n}{h})$
 This window has discontinuities.
Gaussian Kernel: $\phi(u) = \frac{\exp(-\frac{1}{2}\|x\|^2)}{\sqrt{2\pi}}$
 Result in a smoother density model
 $p(x) = \frac{1}{N} \sum_{n=1}^N \frac{1}{(2\pi h^2)^{D/2}} \exp(-\frac{\|x-x_n\|^2}{2h^2})$
 We can chose any other kernel ϕ with $\phi(u) \geq 0 \quad \int \phi(u) du = 1$
K-Nearest Neighbors
 Fix K and find V
 $\hat{p}(x) = \frac{1}{V_k(x)}, v_k(x)$ minimal volume around x containing k neighbors.
Classifier: classify x by the majority of the vote of its k -NN.
1-NN Error Rate the 1-NN error rate P is always $P^* \leq P \leq 2P^*$ where P^* is the error rate of the Bayes rule. \Rightarrow as k goes to infinity kNN becomes optimal
 KNN not optimal if class densities are very different.

Mixture Models
Gaussian Mixture
EM-Algorithm
 Latent Variable: unknown data \rightarrow What cluster generated each sample?
 EM does ML for unknown parameters.

Latent var. $M_{xc} = \begin{cases} 1 & \text{c generated x} \\ 0 & \text{else} \end{cases}$
 $P(\mathcal{X}, M | \theta) = \prod_{x \in \mathcal{X}} \prod_{c=1}^k (\pi_c P(\mathbf{x} | \theta_c))^{M_{xc}}$
E-Step
 $\gamma_{xc} = \mathbb{E}[M_{xc} | \mathcal{X}, \theta^{(j)}] = \frac{P(\mathbf{x} | c, \theta^{(j)}) P(c | \theta^{(j)})}{P(\mathbf{x} | \theta^{(j)})}$

M-Step
 $\mu_c^{(j+1)} = \frac{\sum_{c \in \mathcal{X}} \gamma_{xc} \mathbf{x}}{\sum_{c \in \mathcal{X}} \gamma_{xc}}$
 $(\sigma_c^2)^{(j+1)} = \frac{\sum_{c \in \mathcal{X}} \gamma_{xc} (\mathbf{x} - \mu_c)^2}{\sum_{c \in \mathcal{X}} \gamma_{xc}}$
 $\pi_c^{(j+1)} = \frac{1}{|\mathcal{X}|} \sum_{c \in \mathcal{X}} \gamma_{xc}$

k-Means
 identify clusters of data.
 Given $\mathcal{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$
 Find $c(\cdot)$ and \mathcal{Y} minimizing
 $\mathcal{R}^k m(c, \mathcal{Y}) = \sum_{x \in \mathcal{X}} \|x - \mu_{c(x)}\|^2$ Assign to nearest cluster. Recompute all clusters and repeat. Also called **hard EM**.
 Special case of GMM w. uniform prior and diag. covariance ($\rightarrow 0$).