

# Introduction to Machine Learning summary

Michael Van Huffel, Dino Colombo

August 6, 2022

Introduction to Machine Learning summary created by *michavan@student.ethz.ch* && *dicolomb@student.ethz.ch*

This summary has been written based on the Lecture 252-0220-00 S Introduction to Machine Learning by Prof. A. Krause (Spring 22s). There is no guarantee for completeness and/or correctness regarding the content of this summary. This summary is a corrected, modified and a more completed version of the summary of Yannick Merkli. Use it at your own discretion

**k-CV:**  $k \uparrow, \downarrow$  bias,  $\uparrow$  var,  $R(\text{gen}) \downarrow, \hat{R}(\text{valid}) \uparrow$   
**Bayes:**  $P(X|Y)P(Y) = P(X \cap Y)$   
 $P(X, Y) = P(X|Y)P(Y) = P(Y|X)P(X)$   
 $P(X, Y, Z) = P(X|Y, Z)P(Y|Z)P(Z)$   
 $X, Y \text{ iid}: P(X, Y|Z) = P(X|Z)P(Y|Z)$   
 $\sigma_X^2 = \text{Var}[X] = \mathbb{E}[(X - \mu_X)^2] = \mathbb{E}[X^2] - \mathbb{E}[X]^2$   
 $\ln(x) \leq x - 1, x > 0; \|x\|_2 = \sqrt{x^T x}$

**Jensen ineq:**  $g(\mathbb{E}[X]) \leq \mathbb{E}[g(X)]$ , if  $g \cup \cap$

**Div.**  $KL(p||q) = \mathbb{E}_p[\log(\frac{p(x)}{q(x)})] \neq KL(q||p)$

$JSD(p||q) = \frac{1}{2}[KL(p||\frac{1}{2}(p+q)) + KL(q||\frac{1}{2}(p+q))]$

**Cauchy-Sch.:**  $|\langle x, y \rangle| \leq \|x\| \cdot \|y\|$ , *strong*

**Orth:**  $A^{-1} = A^T, AA^T = A^T A = I, \det = \pm 1$

**Inv:**  $A^{-1} = \frac{1}{\det(A)} \begin{bmatrix} A_{22} & -A_{12} \\ -A_{21} & A_{11} \end{bmatrix}$

**Der:**  $\partial_x(x^T A x) = (A^T + A)x, \partial_x(b^T A x) = A^T b, \partial_x(c^T X^T b) = bc^T, \partial_x a^T X^{-1} b = -X^{-T} ab^T X^{-T}, a^T a = \text{Tr}(aa^T), \text{Tr}(A) = \sum \lambda_i$

**Eigdec:**  $A = Q\Lambda Q^{-1}, \lambda_{1,2}^{2x2} = \frac{\text{Tr}(A) \pm \sqrt{\text{Tr}(A)^2 - 4 \det(A)}}{2}, \mathbf{v}_{1,2} \propto \begin{bmatrix} A_{12} \\ \lambda_{1,2} - A_{11} \end{bmatrix}$

**SVD:**  $X \in \mathbb{R}^{n \times p}, U \in \mathbb{R}^{n \times n}, S \in \mathbb{R}^{n \times p}, V \in \mathbb{R}^{p \times p}$   
 $X = USV^T = \sum_{k=1}^{\text{rank}(X)} \sigma_{k,k} u_k(v_k)^T$   
 $X^T X = V S^T U^T U S V^T = V S^T S V^T = V \Sigma V^T$   
 $\Sigma = \text{diag}(\sigma_i^2); \sigma_i^2 = \lambda_i; \forall \lambda_i \geq 0, (U^T U = V^T V = I)$

**Convex:**  $\forall x_1, x_2, h(\lambda x_1 + (1-\lambda)x_2) \leq \lambda h(x_1) + (1-\lambda)h(x_2)$ , strong if  $h(x) - m/2\|x\|^2 \cup$   
Operations:  $\alpha f + \beta h \cup$  if  $\alpha, \beta \geq 0, f \cup, h \cup$ ;  
 $f(h(x))$  if  $f \cup, h$  affine or  $f$  non-decreasing,  
 $g \cup$ ;  $\max\{f(x), h(x)\}$  if  $f$  convex,  $h$  convex

**Gaussian distribution:**  $\log \prod_i^n f(x_i) = -n \ln \sigma - \frac{n}{2} \ln(2\pi) - \sum_{i=1}^n \frac{1}{2\sigma^2} (x_i - \mu)^2$   
 $f(\mathbf{x}) = \frac{1}{\sqrt{(2\pi)^k |\Sigma|}} \exp(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \Sigma^{-1}(\mathbf{x} - \boldsymbol{\mu}))$

## Regression - S

**Linear Regression**  $f(x) = w^T x, w^* = \min_w \hat{R}$   
 $\hat{R}(w) = \sum_{i=1}^n (y_i - w^T x_i)^2 = \|Xw - y\|_2^2$   
Clf:  $\in \mathcal{O}(nd^2 + d^3)$ :  $w^* = (X^T X)^{-1} X^T y, X \in \mathbb{R}^{n \times d}, w \in \mathbb{R}^{d \times 1}; \nabla_w \hat{R}(w) = -2 \sum_{i=1}^n (y_i - w^T x_i) \cdot x_i = 2X^T(Xw - y) \in \mathcal{O}(nd)$   
**Regularized regression** (*bias*  $\uparrow$ , *variance*  $\downarrow$ )

Ridge/Lasso:  $\min_w \sum_{i=1}^n (y_i - w^T x_i)^2 + \lambda \|w\|_{2,1}^2$

Close f.:  $w^* = (X^T X + \lambda I)^{-1} X^T y = \frac{1}{1+\lambda} w_{OLS}^*$

In general:  $\|w\|_2 \leq \|w\|_1 \leq \sqrt{\|w\|_0} \|w\|_2$

**Gradient descent**  $\mathcal{O}(nd^2 + \tau d^2)$

$\|A\|_{op} = \max\{|\lambda_{\min}|, |\lambda_{\max}|\}$ , SVD for  $X^T X$

$\|w^t - w_{\min}\|_2 \leq \left[ \|I - \eta X^T X\|_{op} \right]^t \|w^0 - w_{\min}\|_2$

**SGD:** 1pt, batch  $\downarrow$ , variance  $\uparrow$ ,  $P(\text{esc. sdl.}) \uparrow$

**k-CV:**  $k \uparrow, \downarrow$  bias,  $\uparrow$  var,  $R(\text{gen}) \downarrow, \hat{R}(\text{valid}) \uparrow$   
**LOOCV:**  $\uparrow \mathcal{O}(\cdot)$ , high var, overfitting,  $= n$ -fold.  
**Classification - S** ( $\ell_{0-1}(f(x), y) = \mathbb{I}_{\{y \neq \text{sign } f(x)\}}$ )  
**Surrogated losses (0-1, non-convex, non-cont.)**  
Hinge loss:  $l_{\text{hinge}}(y\hat{f}(x)) = \max(0, 1 - y\hat{f}(x))$   
Logistic loss:  $l_{\log}(y\hat{f}(x)) = \log(1 + e^{-y\hat{f}(x)})$

CE:  $l_{\text{CE}}(\hat{f}(x), y) = -y \log p(\hat{Y} = 1 | x) - (1 - y) \log(1 - p(\hat{Y} = 1 | x)), Y \in [0, 1]$

Softmax:  $\text{softmax}_{\alpha}(\mathbf{v})_i = \frac{e^{\alpha v_i}}{\sum_{j=1}^K e^{\alpha v_j}}$

## Multi-Class Classification

$\hat{p}_k = \text{softmax}(\hat{f}(x))_k, \hat{f}(x) = (\hat{f}_1(x), \dots, \hat{f}_K(x))$

$\hat{y} = \text{argmax}_{k=1 \dots K} \hat{f}_k(x) = \text{argmax}_{k=1 \dots K} \hat{p}_k$

**OvA:**  $\hat{y}_i = \text{argmax}_{j \in \{1, \dots, c\}} w_j^T x_i$ ; C bin. classif

**OvR:** Train  $\frac{c(c-1)}{2}$  bin. classif., one for each pair (i,j). Class with most positive predictions wins (slower, but no confidence needed)

**Support Vector Machine** (margin =  $1/\|w\|$ )

H:  $\mathcal{L}(w, b, \alpha) = \frac{1}{2} \|w\|_2^2 - \sum_{i=1}^n \alpha_i (y_i (\mathbf{x}_i^T w + b) - 1)$

S:  $\min_w \frac{1}{n} \sum_{i=1}^n (\max\{0, 1 - y_i w^T x_i\} + \lambda \|w\|_2^2)$

$\nabla_w l_i(w) = \begin{cases} -y_i x_i + 2\lambda w & \text{if } y_i w^T x_i < 1 \\ 2\lambda w & \text{if } y_i w^T x_i \geq 1 \end{cases}$

**Metrics: rand. classifier**  $\rightarrow FPR = TPR \forall \tau$   
Accuracy =  $\frac{\# \text{correct predictions}}{\# \text{all predictions}} = \frac{TP + TN}{TP + TN + FP + FN}$

Precis:  $\frac{TP}{\#[\hat{y}=+1]} = \frac{TP}{TP + FP} \rightarrow P(\hat{y} = 1 | \hat{y} = 1)$

Rec  $TPR = \frac{TP}{\#[y=1]} = \frac{TP}{TP + FN} \rightarrow P(\hat{y} = 1 | y = 1)$

$FPR = \frac{FP}{\#[\hat{y}=-1]} = \frac{FP}{TN + FP} \rightarrow P(\hat{y} = 1 | y = -1)$

$FDR = \frac{FP}{\#[\hat{y}=+1]} = 1 - \text{Prec.} \rightarrow P(y = -1 | \hat{y} = 1)$

F1 score =  $\frac{2}{\frac{1}{\text{Precision}} + \frac{1}{\text{Recall}}}, FNR = 1 - TPR$

**Kernels  $k: \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}; x_i^T x_j \rightarrow k(x_i, x_j)$**

**Feature explosion** ( $p = \mathcal{O}\left(\binom{d+m}{m}\right) \approx \frac{(d+m-1) \dots d}{m \dots 1}$ )

Polynomial  $\phi(x)$ ,  $m$  degree polynomial features in dimension  $d$ ,  $n$  features:  $\sim \mathcal{O}(n^2 d^m)$   
 $\uparrow d, \sim \mathcal{O}(n^2 m^d) \uparrow m$ ; with kernel  $k(x, z) = \langle \phi(x), \phi(z) \rangle = (1 + \langle x, z \rangle)^m \sim \mathcal{O}(n^2(d+m))$ .

**Properties of kernel**  $k(x, y) = \phi(x)^T \phi(y)$

**Mercer's Theorem**  $\rightarrow$  infinite feature map  $k$

must be symmetric:  $k(x, y) = k(y, x)$

Kernel matrix must be positive semi-definite.

**Kernel engineering**

$k_1(x, y) + k_2(x, y); k_1(x, y) \cdot k_2(x, y); c \cdot k_1(x, y), c > 0; f(k_1(x, y)), f$  polynomial with positive coefficients or the exponential function.

**Examples of kernels on  $\mathbb{R}^d$**

Linear kernel:  $k(x, y) = x^T y$

Polynomial kernel:  $k(x, y) = (x^T y + 1)^d$

Gaussian kernel:  $k(x, y) = \exp(-\|x - y\|_2^2 / h^2)$

Laplacian kernel:  $k(x, y) = \exp(-\|x - y\|_1 / h)$

$h \uparrow$ , overfit  $\downarrow$ , distr. wider, flatter, B. smoother.  
**Parametric (finite):** LR, linear perceptron.  
**Non-parametric ( $\infty$ ):** ker. perceptron, k-NN  
 $k(x, y) = \min(x, y), |A \cap B|, \frac{1}{\max(\cdot)} = \min(\frac{1}{x}, \frac{1}{y})$   
**Kernelized linear regression**  $\hat{w} = \sum_i \alpha_i \phi(x_i)$   
 $w^* = \min_{\alpha_{1:n}} \sum_{i=1}^n (\sum_{j=1}^n \alpha_j \langle \phi(x_i), \phi(x_j) \rangle - y_i)^2 + \lambda \sum_{i,j} \alpha_i \alpha_j \langle \phi(x_i), \phi(x_j) \rangle = \min_{\alpha} \|\alpha^T K - y\|_2^2 + \lambda \alpha^T K \alpha, \alpha^* = (K + \lambda I)^{-1} y$   
Prediction:  $y = w^{*T} x = \sum_{i=1}^n \alpha_i^* k(x_i, x)$   
**Kernelized Support Vector Machine**  
SVM:  $k_i = [y_1 k(x_i, x_1), \dots, y_n k(x_i, x_n)]$   
 $\min_{\alpha} \sum_{i=1}^n \max\{0, 1 - y_i \alpha^T k_i\} + \lambda \alpha^T K \alpha$   
Prediction:  $y = \text{sign}\left(\sum_{j=1}^n \alpha_j y_j k(x_j, x)\right)$

**k Nearest Neighbor classifier**  
No training, classification during test time:  
 $y = \text{sign}(\sum_{i=1}^n y_i [x_i \text{ among } k \text{ nn of } x])$

**Neural Networks - U/S**  
 $f(x) = \sum_{i=1}^n w_i^{(2)} \phi(\sum_{j=1}^m w_{ij}^{(1)} x_j) = W^{(2)} \phi(W^{(1)} x)$   
 $f(x) = \phi^{(L)}(W^{(L)} \phi^{(L-1)}(W^{(L-1)} \dots (\phi^{(1)}(W^{(1)} x))))$   
**Learning features (loss non convex)**  
Parametr. feat. maps & optimize over params:  
 $w^* = \text{argmin}_{w, \theta} \sum_{i=1}^n l(y_i; f(x))$

One possibility:  $\phi(x, \theta) = \varphi(\theta^T x) = \varphi(z)$

Predict  $y_j = f_j$  reg. /  $y_j = \text{sign}(f_j)$  class.

**Backpropagation**

Error<sub>out</sub>:  $\delta^{(L)} = l'(f) = [l'(f_1), \dots, l'(f_p)]$

Grad<sub>out</sub>:  $\nabla_{W^{(L)}} \ell(W; y, x) = \delta^{(L)} v^{(L-1)T}$

Error<sub>hidden</sub>:  $\delta^{(\ell)} = \phi'(z^{(\ell)}) \odot W^{(\ell+1)T} \delta^{(\ell+1)}$

Grad<sub>hidden</sub>:  $\nabla_{W^{(\ell)}} \ell(W; y, x) = \delta^{(\ell)} v^{(\ell-1)T}$

**Activation functions**

Sigmoid:  $\varphi(z) = \frac{1}{1 + \exp(-z)}; \varphi' = (1 - \varphi(z)) \cdot \varphi(z)$

Tanh<sub>[-1,1]</sub>:  $\varphi(z) = \tanh(z) = \frac{\exp(z) - \exp(-z)}{\exp(z) + \exp(-z)}$

**Initialization (matters) and Var propagation**

Vanish-/explod- gradients: **keep  $\|v_i\|$  bounded**

Dropout: neuron pres  $p \rightarrow \mathbb{E}[I_i] = \mathbb{E}[I_i^2] = p$

$w_i \sim \mathcal{N}(0, \sigma^2) z = \sum_1^d w_i I_i x_i, v = \phi(z), \mathbb{E}[x_i^n] = 1$

$\mathbb{E}[z] = 0, \text{Var}(z) = \sum_1^d \mathbb{E}[x_i^2] \text{Var}(w_i), \mathbb{E}[v^2] \stackrel{!}{=} 1$

ReLU:  $\sigma^2 = \frac{2}{n_{in}}, \text{Tanh: } \sigma^2 = \frac{1}{n_{in}} = \frac{2}{n_{in} + n_{out}}$

**Regularization**

**Ao:** monitor validation set, early stopping; dropout: avoid hidden units memorize training samples; weight decay; batchnorm.

**Batchnorm:**  $\downarrow$  internal cov. shift, enable  $\uparrow \eta_t$ :

1. mini  $\mu_S, \sigma_S^2$ , 2. normalize, 3. scale, shift.

**CNN**  $z_i = \sum_{j=\max(1, i-d+1)}^{\min(i, k)} w_j x_{i-j+1} \mathbb{R}^{d+k-1}$

Filter output size:  $l_{out} = \frac{w_{in} + 2p - f}{s} + 1$

CNN #fts#chan<sub>in</sub>  $\prod$  ker<sub>i</sub>#chan<sub>out</sub> + #chan<sub>out</sub>

ANN #chan<sub>in</sub>  $I_{in}^2$  #chan<sub>out</sub>  $I_{out}^2$  + #chan<sub>out</sub>  $I_{out}^2$

ResN.: skip lay.,  $\uparrow$  deep NN,  $\downarrow$  vanishing grad.

## Clustering - U

**k- means**

$\hat{\mu} = \arg \min_{\mu} \hat{R}(\mu) = \sum_{i=1}^n \min_{j \in \{1, \dots, k\}} \|x_i - \mu_j\|_2^2$

Non-convex, NP-hard, can be kernelized

**Lloyd's heuristic** ( $\mathcal{O}(nkd)$ ,  $d = \text{dim}, n = \#pt$ )

Initialize cluster centers  $\mu^{(0)} = [\mu_1^{(0)}, \dots, \mu_k^{(0)}]$

While still changes in assignments:

$z_i^{(t)} = \min_{j \in \{1 \dots k\}} \|x_i - \mu_j^{(t-1)}\|_2^2; \mu_j^{(t)} = \frac{1}{n_j} \sum_{i: z_i^{(t)} = j} x_i$

**k-mean++** ( $\mathcal{O}(\log k)$ )

Random  $\mu_1^{(0)} := x_i$  for  $i \sim \text{Uniform}(\{1 \dots n\})$

Add centers 2 to  $k$  randomly, proportionally to squared distance to closest selected center

for  $j = 2$  to  $k$ :  $i_j$  sampled with prob.

$\mu_i^{(0)} \leftarrow x_i$  for  $\text{Prob}(i) \propto \min_{l \in \{1 \dots j-1\}} \|x_i - \mu_l^{(0)}\|_2^2$

**Model selection (separated sphere clusters)**

Regularization (favor simple model with few param.); information criterium; elbow method

**Dimension Reduction - U**

**PCA** ( $\approx$  k-m. but  $W$  orth.,  $z_i \notin E_k = e_1 \dots e_k$ )

Given:  $D = \{x_1, \dots, x_n\} \subset \mathbb{R}^d, 1 \leq k \leq d$

$\Sigma_{d \times d} = \frac{1}{n} \sum_{i=1}^n x_i x_i^T, \mu = \frac{1}{n} \sum_{i=1}^n x_i = 0$  !!

$(W, z_1, \dots, z_n) = \text{argmin} \sum_{i=1}^n \|W z_i - x_i\|_2^2$

$W \in \mathbb{R}^{d \times k}$  is orthogonal,  $W^* = (v_1 | \dots | v_k) w / v_i$  evec. of  $\Sigma$  and evals  $\lambda_1 \geq \dots \geq \lambda_d \geq 0$ .

Projections  $z_1, \dots, z_n \in \mathbb{R}^k$  are given by

$z_i = W^T x_i$  with  $\Sigma = \sum_{i=1}^d \lambda_i v_i v_i^T, W^T W = \mathbb{I}$

**Kernel PCA**

$k \geq 1$ , KPCA given by  $\alpha^{(1)}, \dots, \alpha^{(k)} \in \mathbb{R}^n$ ,

$\alpha^{(i)} = \frac{1}{\sqrt{\lambda_i}} v_i$  from:  $K^{n \times n} = \sum_{i=1}^n \lambda_i v_i v_i^T$

Point  $x$  proj. as  $z \in \mathbb{R}^k$ :  $z_i = \sum_{j=1}^n \alpha_j^{(i)} k(x, x_j)$

**Autoencoders**  $f_{en}: \mathbb{R}^d \rightarrow \mathbb{R}^k, f_{de}: \mathbb{R}^k \rightarrow \mathbb{R}^d$

Learn id. fc:  $x \approx f(x; \theta)$ , if  $\varphi(z) = z, AE = PCA$

$f(x; \theta) = f_{de}(f_{en}(x; \theta_{en}); \theta_{de}), W^{(en/de)} = f_{en/de}$

$W^* = \text{argmin}_w \sum_{i=1}^n \|x_i - f(x; \theta)\|_2^2$

**Probability Modeling**

**Risk estimation (iid)**

Estim. risk (unfeasable):  $\tilde{R}(\hat{f}) = \mathbb{E}_x[l(\hat{f}, f^*)]$

Gen risk:  $R(\hat{f}) = \mathbb{E}_{x, y \sim p}[l(\hat{f}, y)] = \tilde{R}(\hat{f}) + \epsilon_{\mathcal{N}}$

$EMSE \propto Bias^2 + Variance + Noise$

$\mathbb{E}_{D, x, y}[(y - \hat{f}_D(x))^2] = \mathbb{E}_x[\mathbb{E}_D[\hat{f}_D(x)] - f^*(x)]^2 + \mathbb{E}_x[\text{Var}_D[\hat{f}_D(x)]] + \mathbb{E}_{x, y}[y - f^*(x)]^2$

Empirical risk:  $\hat{R}_D(w) = \frac{1}{|D|} \sum_{(x, y) \in D} l(\hat{f}, y)$

$\mathbb{E}_{D, val} [\hat{R}_{D, val}(\hat{f}_{D, train})] = R(\hat{f}_{D, train})$

**Risk minimization** ( $f^* = \min_f \mathbb{E}_{x, y \sim p}[l(f, y)]$ )

**Reg:**  $f^* = \mathbb{E}_x[\min_f \mathbb{E}_y[(y - f(x))^2 | X = x]]$

$\int \partial_{\hat{y}}(\hat{y} - y)^2 p(y | x) dy = 0, f^* = \mathbb{E}[y | X = x]$

Predict via:  $\hat{y} = \hat{E}[Y | X = x] = \int y \hat{p}(y | x) dy$

Class:  $f^* = \mathbb{E}_x[\min_{\hat{y}} \mathbb{E}_y[y \mid X = x]] \rightarrow \sum p(y \mid \cdot) [y \neq \hat{y}] = \sum_{y \neq \hat{y}} p(y \mid \cdot) \rightarrow \max_{\hat{y}} p(\hat{y} \mid x)$   
**MLE (min var among all unbiased estim.)**  
 Solve  $\theta^* = \operatorname{argmax}_{\theta} \prod_{i=1}^n \hat{p}(y_i \mid x_i, \theta)$   
 $p_{OLS}(y \mid x, w) = \mathcal{N}(w^\top x, \sigma^2)$   
 $p_C(y \mid x, w) = \operatorname{Ber}(y; \sigma(w^\top x)) = \frac{1}{1 + \exp(-y w^\top x)}$   
**MAP (bias through  $p(\theta \mid x, y) = \frac{p(\theta)p(y, x, \theta)}{p(y \mid x)}$ )**

Gauss prior:  $p(\theta \mid \mu, \beta) = \prod_1^d \frac{1}{\beta \sqrt{2\pi}} e^{-\frac{1}{2}(\frac{\theta_i - \mu_i}{\beta})^2}$   
 Laplace prior:  $p(\theta \mid \mu, b) = \prod_i^d \frac{1}{2b} e^{-\frac{|\theta_i - \mu_i|}{b}}$   
 Ridge:  $p_{\text{gauss}}(\theta \mid 0, \beta), \lambda = \frac{\sigma^2}{\beta^2}, \|\theta\|_2^2 = \sum_i^d \theta_i^2$   
 Lasso:  $p_{\text{laplace}}(\theta \mid 0, b), \lambda = \frac{2\sigma^2}{b}, \|\theta\|_1 = \sum |\theta_i|$   
**MLE for Classification ( $\hat{w}$ )**  
 $\max_w p_C(y_{1:n} \mid x_{1:n}, w) = \min - \sum \log p(y_i \mid x_i, w)$   
 $= \operatorname{argmin}_w \sum_{i=1}^n \log(1 + \exp(-y_i w^\top x_i))$   
**GD for logistic regression**  
 $p(\text{missclas.}) = \hat{p}(Y = -y \mid w, x) = \frac{1}{1 + \exp(y w^\top x)}$   
 Update  $w \leftarrow w + \eta_t y x \hat{p}(Y = +y \mid w, x)$   
 Ridge:  $w \leftarrow w(1 - 2\lambda \eta_t) + \eta_t y x \hat{p}(Y = +y \mid w, x)$   
**Multiclass Logistic Reg. (unique if  $w_c = 0$ )**  
 Estimate  $f_{1-c}(x) = w_1^\top x \dots, w_c^\top x$  for logisitc  
 Predict  $p(y \mid x) = \operatorname{Cat}(y \mid \operatorname{softmax}(f_{1:c}(x)))$   
 $= \frac{\exp(w_y^\top x)}{\sum_i^c \exp(w_i^\top x)} = \frac{1}{2} \exp(w_y^\top x)$  for label  $Y = y$

Overflow problem: let  $s_i = w_i^\top x$ , choose  $s_i^* = \max_{i \in K} s_i$ , subtract  $s_i^*$  from other weights.  
 CE:  $l(y, x; w_1 \dots w_c) = -\log p(Y = y \mid x, w_1 \dots w_c)$   
**Kernelized Logistic Regression**  
 $\hat{\alpha} = \min_{\alpha} \sum_1^n \log(1 + e^{-y_i \alpha^\top K_i}) + \lambda \alpha^\top K \alpha$   
 $\hat{p}(y \mid x, \hat{\alpha}) = \frac{1}{1 + \exp(-y \sum_{j=1}^n \alpha_j k(x_j, x))}$   
**Bayesian decision theory**  
 - Conditional distribution over labels  $p(y \mid x)$   
 - Set of actions  $A$   
 - Cost function  $C: Y \times \mathcal{A} \rightarrow \mathbb{R}$   
 $a^* = \min_{a \in A} \mathbb{E}_y[C(y, a) \mid x] = \sum_y p(y \mid x) C(y, a)$   
**Symm. cost class.**  $C(y, a) = [y \neq a], A = \{\pm 1\}$   
 $a^* = \min_{a \in A} \sum_y p_C(y \mid x) C(y, a) = \operatorname{sign}(w^\top x)$

**Asymmetric cost classification**  $A = \{\pm 1\}$   

$$C(y, a) = \begin{cases} c_{fp}, & \text{if } y = -1 \text{ and } a = +1 \\ c_{fn}, & \text{if } y = +1 \text{ and } a = -1 \\ 0, & \text{otherwise} \end{cases}$$
 $C_+ = \mathbb{E}_y[C(y, 1) \mid x] = p_C(1 \mid x) \cdot 0 + p_C(-1 \mid x) c_{fp}$   
 $C_- = \mathbb{E}[C(y, -1) \mid x] = p_C(1 \mid x) c_{fn} + p_C(-1 \mid x) \cdot 0$   
 Pred.  $+1$  if  $C_+ \leq C_- \rightarrow p(y = 1 \mid x) \geq \frac{c_{fp}}{c_{fp} + c_{fn}}$   
**Doubtful logistic class.**  $A = \{+1, -1, D\}$   

$$C(y, a) = \begin{cases} c_1 \cdot [y \neq a] & \text{if } a \in \{+1, -1\} \\ c_2 & \text{if } a = D \end{cases}$$
  
 $a^* = y$  if  $p_C(y \mid x) \geq 1 - c_2/c_1$ , D otherwise

$p_{OLS}(y \mid x), C(y, a) = (y - a)^2$   
 $A = \mathbb{R}, a^* = \min_a \int C(y, a) p_{OLS}(y \mid x) dy = w^\top x$   
**Asymm. cost reg.**  $\Phi(z) = \int_{-\infty}^z \mathcal{N}(z; 0, 1) dz$   
 $C(y, a) = c_1 \max(y - a, 0) + c_2 \max(a - y, 0)$   
 $A = \mathbb{R} \rightarrow a^* = w^\top x + \sigma \Phi^{-1}(\frac{c_1}{c_1 + c_2})$   
**Active learning (uncertainty sampling, no iid)**  
 Entropy:  $H(X) = -\sum_{i=1}^n p(x_i) \log p(x_i)$   
 Given  $D = D^L \cup D^U$ , initially  $D^L = \emptyset$ .

- Estimate  $p(y \mid x)$  given current  $D^L$   
 - Pick unlabeled  $x_{i_t}^U, i_t \in \min_{x \in D^U} H(p(y \mid x))$   
 - Query label  $y_{i_t}$  and set  $D_L \leftarrow D_L \cup \{(x_{i_t}, y_{i_t})\}$   
**Discriminative vs. Generative Modeling**  
**D:** directly assume  $p(y \mid x)$  and estimate parameters from training dataset.  $\approx$  always S.  
 SVM, NN, KNN, Rand. Forest, Log. regression  $\oplus$  robust,  $\oplus \mathcal{O}(\text{cheaper}) \oplus$  less overfit

**G:** estimate  $p(y)$  and  $p(x \mid y)$  from training data and compute posterior with Bayes's rule:  
 $p(y \mid x) = \frac{p(y)p(x \mid y)}{p(x)}, p(x) = \sum_y p(x, y)$   
 Models: LDA, Naive Bayes, GMM, GANS  
 $\oplus$  outlier detection,  $\oplus$  generate new data,  
 $\oplus$  robust to outliers  $\ominus$  model  $x$  may difficult  
**Outlier detect**  $P(x) = \sum_{y=1}^c P(y) P(x \mid y) \leq \tau$   
**MLE for Categorical( $y \mid \theta$ ),  $\theta = \{p_1, p_0\}$**   
 Let  $p(y = 1) = p_1, p(y = 0) = p_0 = 1 - p_1$   
 $L(y \mid \theta) = \prod_{i=1}^n p_1^{[y_i=1]} p_0^{[y_i=0]} (1 - p_1 - p_0)^{1 - [y_i=1] - [y_i=0]} = p_1^{n_1} p_0^{n_0} (1 - p_1 - p_0)^{1 - n_1 - n_0}$   
 $\nabla_{\theta} \log L(y \mid \theta) \stackrel{!}{=} 0 \Rightarrow p_1 = \frac{n_1}{n_1 + n_0}, p_0 = \frac{n_0}{n_1 + n_0}$

**Gaussian Bayes Classifier (or QDA) - S**  
 MLE class prior:  $P(Y = y) = \hat{p}_y = \frac{\operatorname{Count}(Y=y)}{n}$   
 MLE feature distr:  $P(x \mid y) = \mathcal{N}(x; \hat{\mu}_y, \hat{\Sigma}_y)$   
 $\hat{\mu}_y = \frac{1}{\operatorname{Count}(Y=y)} \sum_{i: y_i=y} x_i \in \mathbb{R}^d$   
 $\hat{\Sigma}_y = \frac{1}{\operatorname{Count}(Y=y)} \sum_{i: y_i=y} (x_i - \hat{\mu}_y)(x_i - \hat{\mu}_y)^T \mathbb{R}^{d \times d}$   
**Naive:** iid,  $\Sigma_y = \operatorname{diag}(\sigma_1^2 \dots \sigma_d^2)$ , predict new  $x$ :  
 $y = \max_y P(y' \mid x) = \max_y P(y') \prod_{i=1}^d P(x_i \mid y')$   
**Decision rule binary classif.**  $c = 2, y \in \{\pm 1\}$   
 $y = \max_y P(y' \mid x) = \operatorname{sign}(f(x)), p_+(x) = \sigma(f(x))$   
 Discriminant fnc:  $f(x) = \log \frac{P(y=1 \mid x)}{P(y=-1 \mid x)} =$   
 $\log \frac{p_+(x)}{1 - p_+(x)} = \log \frac{p}{1 - p} + \frac{1}{2} [\log \frac{|\hat{\Sigma}_-|}{|\hat{\Sigma}_+|} + ((x - \hat{\mu}_-)^T \hat{\Sigma}_-^{-1} (x - \hat{\mu}_-) - ((x - \hat{\mu}_+)^T \hat{\Sigma}_+^{-1} (x - \hat{\mu}_+))]$   
**Naive:**  $p(x \mid y) = \prod_i \mathcal{N}(x_i; \mu_{y,i}, \sigma_{y,i}^2)$   
 If shared  $\sigma_{y,i}^2$ :  $f(x) = w^\top x + w_0 \rightarrow$  log. regress.  
 $w_0 = \log \frac{\hat{p}_+}{1 - \hat{p}_+} + \sum_{i=1}^d \frac{\hat{\mu}_-^2 - i \hat{\mu}_+^2 + i}{2 \hat{\sigma}_i^2}, w_i = \frac{\mu_{+,i} - \mu_{-,i}}{\hat{\sigma}_i^2}$   
**Overcounting:** if  $x_2 = \dots = x_d$  then decision plane stop sum  $w_0/w_i$  at 2.  
**Overconfident:** conditional independence  $\rightarrow$  prediction close to 1 or 0. Fine to pred. most likely class,  $\ominus$  making decision (asymm. loss)

**Linear discriminant analysis (LDA,  $\Sigma_- = \dots = \Sigma_c$ )**  
**Fisher LDA:**  $c = 2, p = 0.5; \hat{\Sigma}_- = \hat{\Sigma}_+ = \hat{\Sigma}$   
 Predict:  $y = \operatorname{sign}(f(x)) = \operatorname{sign}(w^\top x + w_0)$   
 $w = \hat{\Sigma}^{-1}(\hat{\mu}_+ - \hat{\mu}_-); w_0 = \frac{1}{2}(\hat{\mu}_-^\top \hat{\Sigma}^{-1} \hat{\mu}_- - \hat{\mu}_+^\top \hat{\Sigma}^{-1} \hat{\mu}_+)$   
**Categorical Naive Bayes Classifier**  
 MLE class prior:  $P(Y = y) = \hat{p}_y = \frac{\operatorname{Count}(Y=y)}{n}$   
 MLE feature distr:  $P(X_i = x \mid Y = y) = \theta_{x \mid y}^{(i)}$   
 $\hat{\theta}_{x \mid y}^{(i)} = \frac{\operatorname{Count}(X_i=x, Y=y)}{\operatorname{Count}(Y=y)}$ . Prediction as GBC  
**Avoid overfitting**  
 Restrict model class (es only diag cov  $\rightarrow$  few parameters), priors  $\rightarrow$  smaller param. values  
**MAP in Generative Modelling**  
 Use MAP to estimate class prior/feature distrib. ( $\approx 0$  comp. cost). Prior to Bernoulli/Bin:  
 Beta( $\theta; \alpha_+, \alpha_-$ ) =  $\frac{1}{B(\alpha_+, \alpha_-)} \theta^{\alpha_+ - 1} (1 - \theta)^{\alpha_- - 1}$   
 $MLE_{\text{Beta}} = \frac{\alpha_+ - 1}{\alpha_+ + \alpha_- - 2}$   
 Categorical/Multinomial (l)  $\rightarrow$  Dirichlet (prior)  
 Gauss. (l)  $\rightarrow$  Gauss.  $\neg$ -inverse Wishart (prior)  
**Gaussian Mixture Models (GMM) - SSL/U**  
 Sel.  $k$  via CV or  $D_{\text{train}}/D_{\text{val}}$ , maxim. log-like.  
 Can choose  $\tau$  to control estimated  $FPR$   
**Mixture modeling**  $\theta = (\mu, \Sigma, w)$  (iid)  
 Convex comb:  $P(x \mid \theta) = \sum_1^k w_j \mathcal{N}(x; \mu_j, \Sigma_j)$   
 Minimize:  $\theta^* = \min_{\theta} - \sum_i \log \sum_j^k w_j P(x_i \mid \theta_j)$

Non feasible  $\rightarrow$  non-convex, alternatively:  
 $P(x \mid \theta) = \sum_z P(x, z \mid \theta) = \sum_z P(z \mid \theta) P(x \mid z, \theta)$   
 and  $P(z_i = j \mid \theta) = w_j, P(x, z \mid \theta) = w_z \mathcal{N}(x \mid \mu_z, \Sigma_z)$   
 Fitting GMM = training GBC without labels  
**Hard-EM algorithm (most probable  $z_i$ )**  
 Predict class  $z_i$  for each  $x_i$ :  
**E:**  $z_i^{(t)} = \operatorname{argmax}_z P(z \mid x_i, \theta^{(t-1)})$   
 $= \operatorname{argmax}_z P(z \mid \theta^{(t-1)}) P(x_i \mid z, \theta^{(t-1)})$   
 $= \operatorname{argmax}_z w_z^{(t-1)} \mathcal{N}(x_i \mid \mu_z^{(t-1)}, \Sigma_z^{(t-1)})$   
 Now complete:  $D^{(t)} = \{(x_1, z_1^{(t)}) \dots (x_n, z_n^{(t)})\}$   
**M:** MLE as in GBC  $\theta^{(t)} = \operatorname{argmax}_{\theta} P(D^{(t)} \mid \theta)$   
**Issues:** points have label,  $\ominus$  if cluster overlap  
**Lloyd = H-EM**  $w_z = \frac{1}{k}, \Sigma_z = \sigma^2 \mathbb{I}$ , bal class  
 $\rightarrow$  **E:**  $z_i^{(t)} = \operatorname{argmin}_z \|x_i - \mu_z^{(t-1)}\|_2^2$   
 $\rightarrow$  **M:**  $\mu_j^{(t)} = \frac{1}{n_j} \sum_{i: z_i^{(t)}=j} x_i$

**Lloyd h. = S-EM** as above and  $\sigma^2 \rightarrow 0$   
**Soft-EM algorithm (weighting average)**  
 ( $\equiv$  training a GBC with weighted data!)  
 Let iid.  $P(x_{1:n}, z_{1:n} \mid \theta) = \prod_{i=1}^n P(x_i, z_i \mid \theta)$   
**E-step:** calculate expected complete log-like:  
 $Q(\theta; \theta^{(t-1)}) = \mathbb{E}_{z_{1:n}} [\log P(x_{1:n}, z_{1:n} \mid \theta) \mid x_{1:n}, \theta^{(t-1)}]$   
 $= \sum_i \sum_j^k \underbrace{P(z_i = j \mid x_i, \theta^{(t-1)})}_{\gamma_{z_i}^{(t)}(x_i): \text{ex. suf. stat.}} \log \underbrace{P(x_i, z_i = j \mid \theta)}_{w_j \mathcal{N}(x_i \mid \mu_j, \Sigma_j)}$   
 where (Bayes)  $\gamma_j^{(t)}(x_i) = \frac{w_j P(x_i \mid \Sigma_j, \mu_j)}{\sum_{i=1}^n w_i P(x_i \mid \Sigma_i, \mu_i)}$

**M-step:** calculate  $\theta^{(t)} = \operatorname{argmax}_{\theta} Q(\theta; \theta^{(t-1)})$   
 subject to  $\sum_{j=1}^k w_j = 1 \rightarrow$  Lagrange multipliers  
 $\theta^{(t)} = \min_{\theta} \mathcal{L}(\theta, \lambda) = Q(\theta; \theta^{(t-1)}) + \lambda(\sum w_k - 1)$   
 $w_j^{(t)} \leftarrow \frac{1}{n} \sum_{i=1}^n \gamma_j^{(t)}(x_i); \mu_j^{(t)} \leftarrow \frac{\sum_{i=1}^n \gamma_j^{(t)}(x_i) x_i}{\sum_{i=1}^n \gamma_j^{(t)}(x_i)}$   
 $\Sigma_j^{(t)} \leftarrow \frac{\sum_{i=1}^n \gamma_j^{(t)}(x_i) (x_i - \mu_j^{(t)}) (x_i - \mu_j^{(t)})^T}{\sum_{i=1}^n \gamma_j^{(t)}(x_i)} + \nu^2 \mathbb{I}$   
**Avoid degen, calculated via MAP, Wishart prior**  
**Convergence of the Soft-EM algorithm**  
 S-EM monotonically increases likelihood.  
 $\log P(x_{1:n} \mid \theta^{(t)}) \geq \log P(x_{1:n} \mid \theta^{(t-1)})$   
 For Gaussian mixtures, guaranteed to conv. local maxim. Quality solution  $\leftarrow$  initialization!  
 Re-run multiple times, chose sol biggest likh.  
**Semi-S case:**  $\{(x_1, z_1) \dots (x_l, z_l)\}$  labels given  
 For  $i > l \rightarrow \gamma_j^{(t)}(x_i, y_i)$  (usual).

For  $i \leq l \rightarrow \gamma_j^{(t)}(x_i, y_i) = [j = z_i]$   
**Gaussian-Mixture Bayes classifiers**  
 Labeled  $D = \{(x_i, y_i)\}, i = 1 \dots n, m$  labels  
 Estimate class prior  $P(y)$   
 $P(x \mid y) = \sum_{j=1}^{k_y} w_j^{(y)} \mathcal{N}(x; \mu_j^{(y)}, \Sigma_j^{(y)})$   
 $P(y \mid x) = \frac{1}{P(x)} p(y) \sum_{j=1}^{k_y} w_j^{(y)} \mathcal{N}(x; \mu_j^{(y)}, \Sigma_j^{(y)})$   
**Generative adversarial networks (GANS)**  
 Given sample of unlabeled points  $x_1 \dots x_n$   
**Goal:** learn model  $\mathbf{X} = G(\mathbf{Z}; \mathbf{w})$  with  $Z$  simple distrib. (low dim Gauss.),  $G = NN$   
 Likelihood hard, optimize  $w$  to make sample difficult to distinguish from data sample.  
**Training as a "game" betw. 2 NN**  
 1. **Gen.**  $G$  tries produce  $G: z \mapsto G(z; w_G)$   
 2. **Discr.**  $D$  find fake  $D: x \mapsto D(x; w_D) \in [0, 1]$

$$D(x) = D(x, w_D) = \begin{cases} \approx 1 & \text{if } x \text{ real} \\ \approx 0 & \text{if } x \text{ fake} \end{cases}$$
  
 If  $G, D$  enough capacity, guaranteed to conv. to  
**Saddle:**  $\min_{w_G} \max_{w_D} M(w_G, w_D), M(w_G, w_D) =$   
 $\underbrace{\mathbb{E}_{x \sim \text{dat}} \log D(x; w_D)}_{\text{real images}} + \underbrace{\mathbb{E}_{z \sim \mathcal{N}} \log [1 - D(G(z; w_G); w_D)]}_{\text{fake images}}$

**Simultaneous gradient descent:**  
 $w_G^{(t+1)} = w_G^{(t)} - \eta_t \nabla_{w_G} M(w_G, w_D^{(t)})$   
 $w_D^{(t+1)} = w_D^{(t)} + \eta_t \nabla_{w_D} M(w_G^{(t)}, w_D)$   
**Optimal discriminator**  
 For a fixed generator  $G(z; w_G) = x$   
 the optimal discriminator  $D(x)$  is:  
 $P(x \text{ from data}) = D_G^*(x) = \frac{p_{\text{data}}(x)}{p_G(x) + p_{\text{data}}(x)}$   
 $R(G) = \max_D M(G, D) = -\log 4 + 2 JSD(p_{\text{dat}} \parallel p_G)$   
**Duality Gap**  
 $DG(w_G, w_D) = \max_{w_D} M(w_G, w_D') - \min_{w_G} M(w_G', w_D)$   
 $DG = 0$   $w_G, w_D$  pure equilibrium, else  $DG > 0$   
 If  $G$  and  $D$  sufficient capacity, DG upper-bounds the JSD betw. data distrib. and gen.