

Disclaimer

This document is an exam summary that follows the slides of the *Introduction to Machine Learning* lecture at ETH Zurich. The contribution to this is a short summary that includes the most important concepts, formulas and algorithms. This summary was created during the spring semester 2018 by Yannik Merkli and adapted in 2024 by Nicolas Wehrli. Due to updates to the syllabus content, some material may no longer be relevant for future versions of the lecture. This work is published as CC BY-NC-SA.



I do not guarantee correctness or completeness, nor is this document endorsed by the lecturers. Feel free to point out any erratas. For the full \LaTeX source code, consider <https://github.com/nwehrli/ymerkli-eth-summaries>.

Orth: $A: \det(A) \in \{+1, -1\}, AA^T = A^T A = I$
 $trace(ABC) = trace(BCA) = trace(CAB)$
 $trace(A) = \sum \lambda_i(A); \begin{bmatrix} a & b \\ c & d \end{bmatrix}^{-1} = \frac{1}{ad-bc} \begin{bmatrix} d & -b \\ -c & a \end{bmatrix}$
 $A = \sum_{k=1}^{rk(A)} \sigma_{k,k} u_k(v_k)^T, A^\dagger = US'V^T; \sigma'_{k,k} = \frac{1}{\sigma_{k,k}}$
Deriv: $\frac{\partial}{\partial x} b^T x = \frac{\partial}{\partial x} x^T b = b^T, \frac{\partial}{\partial x} \|x\|_2^2 = 2x^T, \frac{\partial}{\partial x} \|x - a\|_2 = \frac{(x-a)^T}{\|x-a\|_2}, \frac{\partial}{\partial x} (x^T A x) = x^T (A^T + A), \frac{\partial}{\partial x} (b^T A x) = A^T b, \nabla_X (c^T X b) = c b^T, \nabla_X (c^T X^T b) = b c^T; \|A\|_{op} = \sup_{\|x\|_2=1} \|A x\|_2$
convex $\iff f(\lambda x + (1-\lambda)y) \leq \lambda f(x) + (1-\lambda)f(y); f(y) \geq f(x) + \langle \nabla f(x), y-x \rangle; D^2 f(x) \geq 0$
 $\alpha f + \beta g$ c.; $\max(f, g)$ c. if f, g c., $\alpha, \beta \geq 0$
 $f \circ g = f(g(x))$ c. if f c., g a. $\forall f$ c., **non-dec.**, g c.
 $p_{\mu, \Sigma}(x) = \frac{1}{\sqrt{(2\pi)^p \det(\Sigma)}} \exp(-\frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu))$
 $X \sim \mathcal{N}(\mu, \Sigma) \implies AX + b \sim \mathcal{N}(A\mu + b, A\Sigma A^T)$
Jensen ineq: $g(E[X]) \leq E[g(X)], g$ convex

Regression
Linear Regression $f(x) = w^T x; X \in \mathbb{R}^{n \times d}$
 $L(w) = \|Xw - y\|_2^2; X^T X \hat{w} = X^T y$
 $d \leq n: \hat{w} = (X^T X)^{-1} X^T y$ if $rk(X) = d$
 $n < d: \hat{w} = (X^T X)^\dagger X^T y; rk(X) = n \implies \|\hat{w}\|_2 \min.$
 $\nabla_w L(w) = 2X^T(Xw - y)$
Gradient Descent
1. Start arbitrary $w_0 \in \mathbb{R}$
2. Do $w_{t+1} = w_t - \eta \nabla L(w_t)$ until $\|w^t - w^{t-1}\|_2 \leq \epsilon$
GD conv. to \hat{w} if $rk(X^T X) = d, \eta < \frac{2}{\lambda_{\max}(X^T X)}$
 $\|w^{t+1} - \hat{w}\| \leq \|I - \eta X^T X\|_{op} \|w^t - \hat{w}\|_2 \leq \rho^{t+1} \|w^0 - \hat{w}\|_2; \eta_{opt} = \frac{2}{\lambda_{\max} + \lambda_{\min}}; \rho_{min} = 1 - \eta_{opt} \lambda_{min} = \frac{\kappa-1}{\kappa+1}$
minibatch SGD: $\nabla L_S(w)$ on random $S \subset D$ every iter.; $|S| = 1$ SGD; $E_S(\nabla L_S(w)) = \nabla L(w)$
strictly c. \implies stationary point is unique g. min.; **strongly c.** \implies unique g. min. exists
Errors
exp. estim. err.: $E_X(\ell(f(X), f^*(X)))$; $y = f^*(x) + \epsilon$
generaliz. err.: $L(f; \mathbb{P}_{X,Y}) = E_{X,Y}(\ell(f(X), Y))$
 $L(\hat{f}; \mathbb{P}_{X,Y}) = E_X((\hat{f}(X) - f^*(X))^2) + \sigma^2$ (sq. loss)
 $L(\hat{f}_D; \mathcal{D}_{test}) = \frac{1}{|\mathcal{D}_{test}|} \sum_{(x,y) \in \mathcal{D}_{test}} \ell(\hat{f}_D(x), y)$ estim.
generaliz. err.; g. err. + const = exp. estim err.
k-fold CV: $\uparrow k \implies \hat{f}_{M_i, D'} \approx \hat{f}_{M_i, D_{use}}, CV_k(M_i) \approx L(\hat{f}_{M_i, D_{use}}; \mathbb{P}_{X,Y})$; extreme: LOOCV

Bias-Variance Tradeoff
 $\text{Bias}_{\mathcal{D}}^2(\hat{f}_D, x) := (E_D(\hat{f}_D(x)) - f^*(x))^2$
 $\text{Bias}_{\mathcal{D}}^2(\hat{f}_D) := E_X(\text{Bias}_{\mathcal{D}}^2(\hat{f}_D, X)); \text{Var}_{\mathcal{D}}(\hat{f}_D) := E_X(\text{Var}_{\mathcal{D}}(\hat{f}_D(X))); E_D(L(\hat{f}_D; \mathbb{P}_{X,Y})) = \text{Var}_{\mathcal{D}}(\hat{f}_D) + \text{Bias}_{\mathcal{D}}^2(\hat{f}_D) + \sigma^2$

Edge closed form: $\hat{w} = (X^T X + \lambda I)^{-1} X^T y$
Classification
Perceptron $y = \text{sign}(f(x)) = \text{sign}(w^T x)$
Perceptron loss is convex and not differentiable, but gradient is informative.
 $l_p(w; y_i, x_i) = \max\{0, -y_i w^T x_i\}$
 $w^* = \arg\min_w \sum_{i=1}^n l_p(w; y_i, x_i)$
 $\nabla_w l_p(w; y_i, x_i) = (-y_i x_i) 1[y_i w^T x_i < 0]$
Stochastic Gradient Descent (SGD)
1. Start at an arbitrary $w_0 \in \mathbb{R}^d$
2. For $t = 1, 2, \dots$ do:
Pick data point $(x', y') \in_{u.a.r.} D$
 $w_{t+1} = w_t - \eta_t \nabla_w l(w; x', y')$
Perceptron Alg: SGD with Perceptron loss
Support Vector Machine
Hinge loss: $l_H(w; x, y) = \max\{0, 1 - y w^T x\}$
Goal: Max. a “band” around the separator.
 $w^* = \arg\min_w \frac{1}{n} \sum_{i=1}^n (\max\{0, 1 - y_i w^T x_i\} + \lambda \|w\|_2^2)$
 $g_i(w) = \max\{0, 1 - y_i w^T x_i\} + \lambda \|w\|_2^2$
 $\nabla_w g_i(w) = \begin{cases} -y_i x_i + 2\lambda w & , \text{ if } y_i w^T x_i < 1 \\ 2\lambda w & , \text{ if } y_i w^T x_i \geq 1 \end{cases}$
Multi-Class Classification
Confidence \rightarrow Distance from Decision Bound.
 $y = \arg\min_{i \in \{1, \dots, c\}} f_i(x), f_i(x) = \bar{w}_i^T x, \bar{w}_i = \frac{w_i}{\|w_i\|_2}$
OvA: $\hat{y}_i = \arg\max_{j \in \{1, \dots, c\}} w_j^T x_i$; C bin. classif
OvO: Train $\frac{c(c-1)}{2}$ bin. classif., one for each pair (i,j). Voting \rightarrow class with most positive predictions wins (slower, but no confidence needed)

Kernels $k: \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}; x_i^T x_j \rightarrow k(x_i, x_j)$
Reformulating the perceptron
Ansatz: $w = \sum_{j=1}^n \alpha_j y_j x_j$
 $w^* = \min_{w \in \mathbb{R}^d} \sum_{i=1}^n \max[0, -y_i w^T x_i]$
 $\iff \alpha^* = \min_{\alpha_{1:n}} \sum_{i=1}^n \max[0, -\sum_{j=1}^n \alpha_j y_j y_i x_i^T x_j]$
Kernelized Perceptron
1. Initialize $\alpha_1 = \dots = \alpha_n = 0$
2. For $t = 1, 2, \dots$ do
Pick data $(x_i, y_i) \in_{u.a.r.} D$
Predict $\hat{y} = \text{sign}(\sum_{j=1}^n \alpha_j y_j k(x_j, x_i))$
If $\hat{y} \neq y_i$ set $\alpha_i^{(t)} = \alpha_i^{(t-1)} + \eta |t|$ else: $\alpha_i^{(t)} = \alpha_i^{(t-1)}$
Predict new point x : $\hat{y} = \text{sign}(\sum_{j=1}^n \alpha_j y_j k(x_j, x))$
Perceptron and SVM
Perceptron: $\min_{\alpha} \sum_{i=1}^n \max\{0, -y_i \alpha^T k_i\}$
SVM: $k_i = [y_1 k(x_i, x_1), \dots, y_n k(x_i, x_n)]$:
 $\min_{\alpha} \sum_{i=1}^n \max\{0, 1 - y_i \alpha^T k_i\} + \lambda \alpha^T D_y K D_y \alpha$
Prediction: $y = \text{sign}(\sum_{j=1}^n \alpha_j y_j k(x_j, x))$

Properties of kernel $k(x, y) = \phi(x)^T \phi(y)$
k must be symmetric: $k(x, y) = k(y, x)$
Kernel matrix must be positive semi-definite.
positive semi-definite matrices
 $M \in \mathbb{R}^{n \times n}$ is psd $\iff \forall x \in \mathbb{R}^n: x^T M x \geq 0 \iff$ all eigenvalues of M are positive: $\lambda_i \geq 0$
k Nearest Neighbor classifier
 $y = \text{sign}(\sum_{i=1}^n y_i [x_i \text{ among } k \text{ nn of } x])$
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: $k(x, y) = \exp(-\|x - y\|_1 / h)$
Kernel engineering
 $k_1(x, y) + k_2(x, y); k_1(x, y) \cdot k_2(x, y); c \cdot k_1(x, y), c > 0; f(k_1(x, y))$, where f is a polynomial with positive coefficients or the exponential function
Kernelized linear regression
Ansatz: $w^* = \sum_i \alpha_i x_i$
Parametric: $w^* = \arg\min_w \sum_i (w^T x_i - y_i)^2 + \lambda \|w\|_2^2$
 $= \arg\min_{\alpha} \alpha^T K \alpha - y^T \alpha + \lambda \alpha^T K \alpha, \alpha^* = (K + \lambda I)^{-1} y$
Prediction: $y = w^{*T} x = \sum_{i=1}^n \alpha_i^* k(x_i, x)$

Imbalance
Cost Sensitive Classification
Replace loss by: $l_{CS}(w; x, y) = c_y l(w; x, y)$
Metrics (convention: positive = rare)
Accuracy = $\frac{\# \text{correct predictions}}{\# \text{all predictions}} = \frac{TP + TN}{TP + TN + FP + FN}$,
Precision = $\frac{\# \text{correct}' + \text{'predictions}}{\# \text{all}' + \text{'predictions}} = \frac{TP}{TP + FP}$
Recall = TPR = $\frac{TP}{TP + FN} = \frac{TP}{n_+}$, FPR = $\frac{FP}{TN + FP} = \frac{FP}{n_-}$
F1 score = $\frac{2TP}{2TP + FP + FN} = \frac{2}{\frac{1}{\text{Precision}} + \frac{1}{\text{Recall}}}$

Multi-class Hinge Loss
 $l_{MC-H}(w^{(1)}, \dots, w^{(c)}; x, y) = \max_{j \in \{[1, \dots, c] \setminus y\}} (0, 1 + \max_{i \in y} w^{(i)T} x - w^{(y)T} x)$
Neural Networks
 $F(x) = \sum_{i=1}^k w_i^{(2)} \phi(\sum_{j=1}^m w_{ij}^{(1)} x_j) = W^{(2)} \phi(W^{(1)} x)$
 $F(x) = \phi^{(L)}(W^{(L)} \phi^{(L-1)}(W^{(L-1)} \dots (\phi^{(1)}(W^{(1)} x) \dots)))$
Learning features
Parametr. feat. maps & optimize over params:
 $w^* = \arg\min_{w, \theta} \sum_{i=1}^n l(y_i; \sum_{j=1}^m w_j \phi(x_i, \theta_j))$
One possibility: $\phi(x, \theta) = \varphi(\theta^T x) = \varphi(z)$
Activation functions
Sigmoid: $\varphi(z) = \frac{1}{1 + \exp(-z)}$; $\varphi'(z) = (1 - \varphi(z)) \cdot \varphi(z)$
Tanh $_{[-1,1]}$: $\varphi(z) = \tanh(z) = \frac{\exp(z) - \exp(-z)}{\exp(z) + \exp(-z)}$
ReLU: $\varphi(z) = \max(z, 0)$

Forward propagation
For each unit j on input layer, set value $v_j = x_j$
For each layer $l = 1 : L - 1$: For each unit j

on layer l set its value $v_j = \varphi(\sum_{i \in \text{Layer}_{l-1}} w_{j,i} v_i)$
For each unit j on output layer, set its value
 $f_j = \sum_{i \in \text{Layer}_{L-1}} w_{j,i} v_i$ resp. $\vec{f} = W^{(L)} v^{(L-1)}$
Predict $y_j = f_j$ for reg. / $y_j = \text{sign}(f_j)$ for class.
Backpropagation
For each unit j on the output layer L :
- Compute error signal: $\delta_j^{(L)} = \ell_j'(f_j)$
- For each unit i on layer $L - 1$: $\frac{\partial l}{\partial w_{j,i}^{(L)}} = \delta_j^{(L)} v_i^{(L-1)}$
For each unit j on hidden layer $l = \{L - 1, \dots, 1\}$:
- Error sig: $\delta_j^{(l)} = \varphi'(z_j^{(l)}) \sum_{i \in \text{Layer}_{l+1}} w_{i,j}^{(l+1)} \delta_i^{(l+1)}$
- For each unit i on layer $l - 1$: $\frac{\partial l}{\partial w_{j,i}^{(l)}} = \delta_j^{(l)} v_i^{(l-1)}$

Learning with momentum
 $a \leftarrow m \cdot a + \eta_t \nabla_W l(W; y, x); W \leftarrow W - a$
Clustering
k-mean
Loss: $\hat{R}(\mu) = \hat{R}(\mu_1, \dots, \mu_k) = \sum_{i=1}^n \min_{j \in \{1, \dots, k\}} \|x_i - \mu_j\|_2^2$
 $\hat{\mu} = \arg\min_{\mu} \hat{R}(\mu)$; non-convex, $\mathcal{O}(NP)$

Algorithm (Lloyd’s heuristic):
Initialize cluster centers $\mu^{(0)} = [\mu_1^{(0)}, \dots, \mu_k^{(0)}]$
While still changes in assignments:
 $z_i = \arg\min_{j \in \{1, \dots, k\}} \|x_i - \mu_j^{(t-1)}\|_2^2; \mu_j^{(t)} = \frac{1}{n_j} \sum_{i: z_i=j} x_i$
k-mean++:
- Start with random data point as center
- Add centers 2 to k randomly, proportionally to squared distance to closest selected center for $j = 2$ to k : i_j sampled with prob.
 $P(i_j = i) = \frac{1}{z} \min_{1 \leq l < j} \|x_i - \mu_l\|_2^2; \mu_j \leftarrow x_{i_j}$

Dimension Reduction
Principal component analysis (PCA)
Given: $D = \{x_1, \dots, x_n\} \subset \mathbb{R}^d, 1 \leq k \leq d$
 $\Sigma_{d \times d} = \frac{1}{n} \sum_{i=1}^n x_i x_i^T, \mu = \frac{1}{n} \sum_{i=1}^n x_i = 0 !!$
Sol.: $(W, z_1, \dots, z_n) = \arg\min_{\sum_{i=1}^n \|W z_i - x_i\|_2^2}$, where $W \in \mathbb{R}^{d \times k}$ is orthogonal, $W^* = (v_1 | \dots | v_k)$ w/ v_i evec. of Σ and evals $\lambda_1 \geq \dots \geq \lambda_d \geq 0$.
Projections $z_1, \dots, z_n \in \mathbb{R}^k$ are given by
 $z_i = W^T x_i$ where $\Sigma = \sum_{i=1}^d \lambda_i v_i v_i^T$,
Kernel PCA
For general $k \geq 1$, the Kernel PC are given by $\alpha^{(1)}, \dots, \alpha^{(k)} \in \mathbb{R}^n$, where $\alpha^{(i)} = \frac{1}{\sqrt{\lambda_i}} v_i$ is obtained from: $K = \sum_{i=1}^n \lambda_i v_i v_i^T, \lambda_1 \geq \dots \geq \lambda_d \geq 0$
Point x projected as $z \in \mathbb{R}^k: z_i = \sum_{j=1}^n \alpha_j^{(i)} k(x, x_j)$

Useful coders: $f_1 : \mathbb{R}^d \rightarrow \mathbb{R}^k, f_2 : \mathbb{R}^k \rightarrow \mathbb{R}^d$
 Try to learn identity function: $x \approx f(x; \theta)$
 $f(x; \theta) = f_2(f_1(x; \theta_1); \theta_2); f_1 : \text{en-}, f_2 : \text{decoder}$
 d input, d output units, 1 layer w/ $k < d$ units
 $W^* = \text{argmin}_w \sum_{i=1}^n \|x_i - W^{(2)} \varphi(W^{(1)} x^{(i)})\|_2^2$
 $\varphi(z) = z : NNA = PCA, w^{(1)} = PCA(x) = w^{(2)T}$
Probability Modeling

Assumption: Data set is generated iid
 Find $h : X \rightarrow Y$ that minimizes pred. error
 $R(h) = \int P(x, y) l(y; h(x)) dx dy = \mathbb{E}_{x, y} [l(y; h(x))]$
 $h^*(x) = \mathbb{E}[Y|X=x]$ for $R(h) = \mathbb{E}_{x, y} [(y - h(x))^2]$
 Prediction: $\hat{y} = \hat{\mathbb{E}}[Y|X=x] = \int \hat{P}(y|X=x) y dy$

Maximum Likelihood Estimation (MLE)
 Choose a particular parametric form $\hat{P}(Y|X, \theta)$
 $\theta^* = \text{amax}_{\theta} \hat{P}(y_{1:n}|x_{1:n}, \theta) \stackrel{\text{iid}}{=} \text{amax}_{\theta} \prod_{i=1}^n \hat{P}(y_i|x_i, \theta)$
 $= \text{amin}_{\theta} - \sum_{i=1}^n \log \hat{P}(y_i|x_i, \theta)$

Ex: $y_i \sim \mathcal{N}(w^T x_i, \sigma^2) : w^* = \text{amin}_w \sum_i (y_i - w^T x_i)^2$
Bias/Variance/Noise
 Prediction error = $\text{Bias}^2 + \text{Variance} + \text{Noise}$

Maximum a posteriori estimate (MAP)
 Introduce bias by expressing assumption through a Bayesian prior $w_i \sim \mathcal{N}(0, \beta^2)$
 Bayes: $P(w|x, y) = \frac{P(w|x)P(y|x, w)}{P(y|x)} = \frac{P(w)P(y|x, w)}{P(y|x)}$
 assume w indep. of x: $\text{argmax}_w P(w|x, y) = \text{argmin}_w -\log P(w) - \log P(y|x, w) + \text{const.}$
 $= \text{argmin}_w \lambda \|w\|_2^2 + \sum_{i=1}^n (y_i - w^T x_i)^2, \lambda = \frac{\sigma^2}{\beta^2}$
 (= $\text{argmax}_w P(w) \prod_i P(y_i|x_i, w)$, assuming noise w)

$P(y|x, w)$ iid Gaussian, prior $P(w)$ Gaussian)
Logistic regression
 Assume iid Bernoulli noise instead of Gauss.
 $P(y|x, w) = \text{Ber}(y; \sigma(w^T x)) = \frac{1}{1 + \exp(-y w^T x)}$

$l_{\text{logistic}}(w; x_i, y_i) = \log(1 + \exp(-y_i w^T x_i))$
 $\nabla_w l(w) = \frac{(-y_i x_i)}{1 + \exp(+y_i w^T x_i)} = P(Y = -y|x, w)(-y_i x_i)$

Example: MLE for logistic regression
 $\text{argmax}_w P(y_{1:n}|w, x_{1:n})$
 $= \text{argmin}_w - \sum_{i=1}^n \log P(y_i|w, x_i)$
 $= \text{argmin}_w \sum_{i=1}^n \log(1 + \exp(-y_i w^T x_i))$
 $\hat{R}(w) = \sum_{i=1}^n \log(1 + \exp(-y_i w^T x_i))$ (neg log l. f.)
Logistic regression and regularization
 $\min_w \sum_{i=1}^n \log(1 + \exp(-y_i w^T x_i)) + \lambda (\|w\|_1 \text{ or } \|w\|_2^2)$

SGD for logistic regression
 Update $w \leftarrow w + \eta_t y x \hat{P}(Y = -y|w, x)$
L2 regularized logistic regression:
 Update $w \leftarrow w(1 - 2\lambda \eta_t) + \eta_t y x \hat{P}(Y = -y|w, x)$
Multiclass Logistic Regression
 $P(Y = i|x, w_1, \dots, w_c) = \exp(w_i^T x) / \sum_j \exp(w_j^T x)$

Bayesian decision theory
 - Conditional distribution over labels $P(y|x)$
 - Set of actions \mathcal{A} - Cost function $C : Y \times \mathcal{A} \rightarrow \mathbb{R}$
 Pick action that minimizes the expected cost:
 $a^* = \text{argmin}_{a \in \mathcal{A}} \mathbb{E}_y [C(y, a)|x] = \sum_y P(y|x) C(y, a)$
 $\mathbb{E}_y [C(y, +)|x] = P(-|x) C(-, +);$
 $\mathbb{E}_y [C(y, -)|x] = P(+|x) C(+, -);$
 $\mathbb{E}_y [C(y, D)|x] = P(+|x) C_{d+} + P(-|x) C_{d-}$

Optimal decision for logistic regression
 $a^* = \text{argmin}_y \hat{P}(y|x) = \text{sign}(w^T x)$
Doubtful logistic regression
 Est. cond. distr.: $\hat{P}(y|x) = \text{Ber}(y; \sigma(\hat{w}^T x))$
 Action set: $\mathcal{A} = \{+1, -1, D\}$; Cost function:

$C(y, a) = \begin{cases} [y \neq a] & \text{if } a \in \{+1, -1\} \\ c & \text{if } a = D \end{cases}$
 $\rightarrow a^* = y$ if $\hat{P}(y|x) \geq 1 - c$, D otherwise

Linear regression
 Est. cond. distr.: $\hat{P}(y|x, w) = \mathcal{N}(y; w^T x, \sigma^2)$
 $\mathcal{A} = \mathbb{R}; C(y, a) = (y - a)^2$
 $\rightarrow a^* = \mathbb{E}_y [y|x] = \int \hat{P}(y|x) dy = \hat{w}^T x$

Asymmetric cost for regression
 Est. cond. distr.: $\hat{P}(y|x) = \mathcal{N}(\hat{y}; \hat{w}^T x, \sigma^2)$
 $\mathcal{A} = \mathbb{R}; C(y, a) = c_1 \max(y - a, 0) + c_2 \max(a - y, 0)$
 $\rightarrow a^* = \hat{w}^T x + \sigma \Phi^{-1}(\frac{c_1 - c_2}{c_1 + c_2}), \Phi$: Gaussian CDF

Discriminative vs. Generative Modeling
 Discriminative models: aim to estimate $P(y|x)$
 G. m.: aim to estimate joint distribution $P(y, x)$
 Typical approach to generative modeling:
 - Estimate prior on labels $P(y)$
 - Estimate cond. distr. $P(x|y)$ for each class y
 - Obtain predictive distr. using Bayes' rule:

$P(y|x) = \frac{P(y)P(x|y)}{P(x)} = \frac{P(x, y)}{P(x)}, P(x) = \sum_y P(x, y)$

Example MLE for P(y)
 Want: $P(Y = 1) = p, P(Y = -1) = 1 - p$
 Given: $D = \{(x_1, y_1), \dots, (x_n, y_n)\}$
 $P(D|p) = \prod_{i=1}^n p^{1[y_i=+1]} (1-p)^{1[y_i=-1]}$
 $= p^{n_+} (1-p)^{n_-}$, where $n_+ = \#$ of $y = +1$

$\frac{\partial}{\partial p} \log P(D|p) = n_+ \frac{1}{p} - n_- \frac{1}{1-p} \stackrel{!}{=} 0 \Rightarrow p = \frac{n_+}{n_+ + n_-}$

Example MLE for P= (x|y)
 Assume: $P(X = x_i|y) = \mathcal{N}(x_i; \mu_{i, y}, \sigma_{i, y}^2)$
 Given: $D, D_{x_i|y} = \{x, \text{s.t. } x_{j, i} = x, y_j = y\}$
 Thus MLE yields: $\hat{\mu}_{i, y} = \frac{1}{n_y} \sum_{x \in D_{x_i|y}} x;$

$\hat{\sigma}_{i, y}^2 = \frac{1}{n_y} \sum_{x \in D_{x_i|y}} (x - \hat{\mu}_{i, y})^2$
Deriving decision rule
 $P(y|x) = \frac{1}{2} P(y) P(x|y), Z = \sum_y P(y) P(x|y)$
 $y = \text{argmax}_y P(y'|x) = \text{argmax}_y P(y') \prod_{i=1}^d P(x_i|y')$
 $= \text{argmax}_y \log P(y') + \sum_{i=1}^d \log P(x_i|y')$

Gaussian Naive Bayes Classifier
 Indep. feat. giv. Y: $P(X_1, \dots, X_n|Y) = \prod_{i=1}^d P(X_i|Y)$
 MLE for class prior: $\hat{P}(Y = y) = \hat{p}_y = \frac{\text{Count}(Y=y)}{n}$
 MLE for feature distr.: $\hat{P}(x_i|y_i) = \mathcal{N}(x_i; \hat{\mu}_{y, i}, \sigma_{y, i}^2)$
 $\hat{\mu}_{y, i} = \frac{1}{\text{Count}(Y=y)} \sum_{j: y_j=y} x_{j, i}$
 $\sigma_{y, i}^2 = \frac{1}{\text{Count}(Y=y)} \sum_{j: y_j=y} (x_{j, i} - \hat{\mu}_{y, i})^2$
 Prediction given new point x:

$y = \text{argmax}_{y'} \hat{P}(y'|x) = \text{argmax}_{y'} \hat{P}(y') \prod_{i=1}^d \hat{P}(x_i|y')$
Categorical Naive Bayes Classifier

MLE class prior: $\hat{P}(Y = y) = p_y = \frac{\text{Count}(Y=y)}{n}$
 MLE for feature distr.: $\hat{P}(X_i = c|Y = y) = \theta_{c|y}^{(i)}$
 $\theta_{c|y}^{(i)} = \frac{\text{Count}(X_i=c, Y=y)}{\text{Count}(Y=y)}, \text{Pred: } y = \text{amax}_{y'} \hat{P}(y'|x)$
 Discr fnc: $f(x) = \log \frac{P(y=1|x)}{P(y=-1|x)}; p(x) = \frac{1}{1 + \exp(-f(x))}$

Gaussian Bayes Classifier
 MLE for class prior: $\hat{P}(Y = y) = \hat{p}_y = \frac{\text{Count}(Y=y)}{n}$
 MLE for feature distr.: $\hat{P}(x|y) = \mathcal{N}(x; \hat{\mu}_y, \hat{\Sigma}_y)$
 $\hat{\mu}_y = \frac{1}{\text{Count}(Y=y)} \sum_{i: y_i=y} x_i \in \mathbb{R}^d$
 $\hat{\Sigma}_y = \frac{1}{\text{Count}(Y=y)} \sum_{i: y_i=y} (x_i - \hat{\mu}_y)(x_i - \hat{\mu}_y)^T \in \mathbb{R}^{d \times d}$

Fisher's linear discriminant analysis (LDA; c=2)
 Assume: $p = 0.5; \hat{\Sigma}_- = \hat{\Sigma}_+ = \hat{\Sigma}$

discriminant f.: $f(x) = \log \frac{p}{1-p} + \frac{1}{2} [\log \frac{|\hat{\Sigma}_-|}{|\hat{\Sigma}_+|} + ((x - \hat{\mu}_-)^T \hat{\Sigma}_-^{-1} (x - \hat{\mu}_-) - ((x - \hat{\mu}_+)^T \hat{\Sigma}_+^{-1} (x - \hat{\mu}_+))]$
 Predict: $y = \text{sign}(f(x)) = \text{sign}(w^T x + w_0)$
 $w = \hat{\Sigma}^{-1} (\hat{\mu}_+ - \hat{\mu}_-); w_0 = \frac{1}{2} (\hat{\mu}_-^T \hat{\Sigma}^{-1} \hat{\mu}_- - \hat{\mu}_+^T \hat{\Sigma}^{-1} \hat{\mu}_+)$

Outlier Detection
 $P(x) = \sum_{y=1}^c P(y) P(x|y) = \sum_y \hat{p}_y \mathcal{N}(x; \hat{\mu}_y, \hat{\Sigma}_y) \leq \tau$

Latent: Missing Data (Gaussian distr.)
Mixture modeling $P(x|\theta) = P(x|\mu; \Sigma, w)$

1) Model each cluster j as prob. distr. $P(x|\theta_j)$
 2) data iid, lklh.: $P(D|\theta) = \prod_{i=1}^n \sum_{j=1}^k w_j P(x_i|\theta_j)$
 3) θ should minimize neg log-likelihood:
 $\theta^* = \text{amin}_{\theta} L(D; \theta) = \text{amin}_{\theta} - \sum_i \log \sum_j w_j P(x_i|\theta_j)$

Ex: $P(x|\theta) = \sum_i w_i \mathcal{N}(x; \mu_i, \Sigma_i), P(z_i = j) = w_j$
 $\sum w_i = 1, P(z, x) = w_z \mathcal{N}(x|\mu_z, \Sigma_z)$

Gaussian-Mixture Bayes classifiers
 Estimate class prior $P(y)$; Est. cond. distr. for

each class: $P(x|y) = \sum_{j=1}^{k_y} w_j^{(y)} \mathcal{N}(x; \mu_j^{(y)}, \Sigma_j^{(y)})$

$P(y|x) = \frac{1}{P(x)} p(y) \sum_{j=1}^{k_y} w_j^{(y)} \mathcal{N}(x; \mu_j^{(y)}, \Sigma_j^{(y)})$

Hard-EM algorithm
 Initialize parameters $\theta^{(0)}$
 For $t = 1, 2, \dots$: Predict class z_i for each x_i :
 $\mathbf{E}: z_i^{(t)} = \text{argmax}_z P(z|x_i, \theta^{(t-1)}) =$

$= \text{argmax}_z P(z|\theta^{(t-1)}) P(x_i|z, \theta^{(t-1)}) =$
 $= \text{argmax}_z w_z^{(t-1)} \mathcal{N}(x_i|\mu_z^{(t-1)}, \Sigma_z^{(t-1)})$
M: Compute the MLE as for the Gaussian B. class.: $\theta^{(t)} = \text{argmax}_{\theta} P(D^{(t)}|\theta)$
 Special case: fix $w_z = \frac{1}{k}$, spher. cov. $\Sigma_z = \sigma^2 \mathbb{I}$
 \rightarrow k-means: **E:** $z_i^{(t)} = \text{argmin}_z \|x_i - \mu_z^{(t-1)}\|_2^2$
M: $\mu_j^{(t)} = \frac{1}{n_j} \sum_{i: z_i^{(t)}=j} x_i$

Soft-EM algorithm: While not converged
E-step: For each i and j calculate $\gamma_j^{(t)}(x_i)$

$\gamma_j^{(t)}(x_i) = P(Z_i = j|x_i, \theta_t) = \frac{P(x_i|Z_i=j, \theta_t) P(Z_i=j|\theta_t)}{P(x_i; \theta_t)} =$
 $= \frac{w_j P(x|\Sigma_j, \mu_j)}{\sum_l w_l P(x|\Sigma_l, \mu_l)} = \frac{w_j \mathcal{N}(x; \Sigma_j, \mu_j)}{\sum_l w_l \mathcal{N}(x; \Sigma_l, \mu_l)}$

$Q(\theta; \theta^{(t-1)}) = \mathbb{E}_{y_{1:n}} [\log P(x_{1:n}, y_{1:n}|\theta) | x_{1:n}, \theta^{(t-1)}] =$
 $= \mathbb{E}_{y_{1:n}} [\log \prod_{i=1}^n P(x_i, y_i|\theta) | x_{1:n}, \theta^{(t-1)}] =$
 $= \sum_{i=1}^n \mathbb{E}_{y_i} [\log P(x_{1:n}, y_i|\theta) | x_i, \theta^{(t-1)}] =$
 $\sum_{i=1}^n \sum_{j=1}^k P(y_i = j | x_i, \theta^{(t-1)}) \log(P(x_i, y_i = j; \theta))$
 $= \sum_{i=1}^n \sum_{j=1}^k \gamma_j^{(t)}(x_i) \log(P(y_i = j) P(x_i|y_i = j; \theta))$
 If constraint $\sum_{j=1}^m P(y_i = j; \theta) = 1$ (m: #labels):
 $\rightarrow \mathcal{L}(\theta, \lambda) = Q(\theta; \theta^{(t-1)}) + \lambda (\sum_j^m P(y_i = j) - 1)$

M-step: Fit clusters to weighted data points:
 Genrl: $\theta^{(t)} = \text{argmax}_{\theta} Q(\theta; \theta^{(t-1)}), \gamma_j^{(t)}(x_i) \text{ fixed!}$

$w_j^{(t)} \leftarrow \frac{1}{n} \sum_{i=1}^n \gamma_j^{(t)}(x_i); \mu_j^{(t)} \leftarrow \frac{\sum_{i=1}^n \gamma_j^{(t)}(x_i) x_i}{\sum_{i=1}^n \gamma_j^{(t)}(x_i)}$

$\Sigma_j^{(t)} \leftarrow \frac{\sum_{i=1}^n \gamma_j^{(t)}(x_i) (x_i - \mu_j^{(t)}) (x_i - \mu_j^{(t)})^T}{\sum_{i=1}^n \gamma_j^{(t)}(x_i)} \{+ \nu^2 \mathbb{I}\}$

SSL w/ GMMs: labeled p.: $y_i: \gamma_j^{(t)}(x_i) = 1 [j = y_i]$
 unl. p.: $\gamma_j^{(t)}(x_i) = P(Z = j|x_i, \mu^{(t-1)}, \Sigma^{(t-1)}, w^{(t-1)})$
Additions: small variance & high bias \rightarrow too simple model

CNN filter output size: $L = \frac{n+2p-f}{s} + 1$