*Projection*: $\text{proj}_b(a) = \frac{a^\top b}{\|b\|^2} b$.

Bayes. has 3 steps: (1) def. prior, (2) def. likelihood, and (3) Bayes rule to compute posterior.

*Cholesky decomp*: If $\mathbf{A}$ is PD, then $\exists \mathbf{L}$ s.t. $\mathbf{A} = \mathbf{L}^\top \mathbf{L}$.

*Geometric series*: $\sum_{k=0}^\infty ar^k = a/1-r$ if $|r| < 1$.

$1 - x \leq \exp(-x) \implies (1-\epsilon)^n \leq \exp(-n\epsilon)$.

**Information theory:**

$$H(p) = \mathbb{E}_p[-\log p(X)]$$
$$D_{\text{KL}}(p \parallel q) = \mathbb{E}_p[\log p(X)/q(X)]$$
$$H(p,q) = \mathbb{E}_p[-\log q(X)] = H(p) + D_{\text{KL}}(p \parallel q)$$
$$\text{I}(X;Y) = \mathbb{E}[\log p(X,Y)/p(X)p(Y)] = H(X) - H(X \mid Y).$$

**Gaussian:**
$(2\pi)^{-n/2}|\mathbf{\Sigma}|^{-1/2} \exp\left(-\frac{1}{2}(x-\mu)^\top \mathbf{\Sigma}^{-1}(x-\mu)\right)$.

*Conditional*: If $\begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \sim \mathcal{N}\left(\begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}, \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix}\right)$.
Then: $x_2 \mid x_1 = z \sim \mathcal{N}(\bar{\mu}, \bar{\mathbf{\Sigma}})$, where $\bar{\mu} = \mu_2 + \Sigma_{21}\Sigma_{11}^{-1}(z - \mu_1)$ and $\bar{\mathbf{\Sigma}} = \Sigma_{22} - \Sigma_{21}\Sigma_{11}^{-1}\Sigma_{12}$.

*Information theory*:

$$D_{\text{KL}} = \frac{1}{2}\left[\log \frac{|\Sigma_2|}{|\Sigma_1|} - d + \text{tr}(\Sigma_2^{-1}\Sigma_1) + (\mu_2 - \mu_1)^\top \Sigma_2^{-1}(\mu_2 - \mu_1)\right]$$
$$H = \frac{d}{2}\log(2\pi e) + \frac{1}{2}\log|\mathbf{\Sigma}|.$$

### Paradigms of data science

Frequentism (optimize likelihood, MLE):
$\theta^\star \in \text{argmax}_{\theta \in \Theta} \sum_{i=1}^n \log p(x_i \mid \theta)$.
Bayesianism (optimize posterior, MAP):
$\theta^\star \in \text{argmax}_{\theta \in \Theta} \log p(\theta) + \sum_{i=1}^n \log p(x_i \mid \theta)$.
Statistical learning (optimize risk):
$f^\star \in \text{argmax}_{f \in \mathcal{F}} \mathcal{R}(f) \doteq \mathbb{E}_{X,Y}[\ell(Y, f(X))]$
$\hat{f}_n \in \text{argmax}_{f \in \mathcal{F}} \hat{\mathcal{R}}_n(f) \doteq \frac{1}{n}\sum_{i=1}^n \ell(y_i, f(x_i))$.

### Anomaly detection

Objects $\mathcal{X} \subseteq \mathbb{R}^d$ with normal class $\mathcal{N} \subseteq \mathcal{X}$. Construct $\phi : \mathcal{X} \to \{0,1\}$ such that $\phi(x) = \mathbb{1}\{x \notin \mathcal{N}\}$. Anomaly is an "unlikely event" $\Rightarrow$ Fit distribution to $\mathcal{X}$ and score according to $p(x)$.

**PCA:** Proj. $\mathcal{X}$ to low-dim. $\Rightarrow \Pi(\mathcal{N})$ is simpler.

Linearly project $\mathbb{R}^d$ to $\mathbb{R}^{d^-}$ such that maximum variance is preserved. Base case $d^- = 1$: Find $u$ with $\|u\| = 1$ s.t. $x \mapsto u^\top x$. Sample mean and variance of reduced dataset:
$\mathbb{E}[u^\top x] = u^\top \mathbb{E}[x]$ and $\mathbb{V}[u^\top x] = u^\top \text{Cov}(x)u$.
We want maximum variance so we have:
$u^\star \in \text{argmax}_{\|u\|=1} u^\top \text{Cov}(x)u$. Solvable by vanishing Lag. grad. Easy to find that $u^\star$ is eigenvector with maximum eigenvalue. Then project it out ($\mathcal{X}_1 = \{x - \text{proj}_{u_1}(x)\} = \{x - u_1^\top x \cdot u_1\}$) and do the same for next dimension

**GMM:** Lin. proj. onto low-dim. spaces resemble Gaussian dist. $\Rightarrow$ Fit GMM to $\Pi(\mathcal{X})$.

Fit $p(x;\theta) = \sum_{j=1}^k \pi_j \mathcal{N}(x; \mu_j, \Sigma_j)$ to data with EM algorithm. Can derive $\log p(\mathbf{X};\theta) = M(q,\theta) + \mathbb{E}(q,\theta)$, where $M(q,\theta) \doteq \mathbb{E}_q[\log p(\mathbf{X},z;\theta)/q(z)]$ and $E(q,\theta) \doteq \mathbb{E}_q[\log q(z)/p(z|\mathbf{X};\theta)]$. Properties: $\log p(\mathbf{X};\theta) \geq M(q,\theta)$ and $\log p(\mathbf{X};\theta) = M(q^\star,\theta)$ where $q^\star = p(\cdot \mid \mathbf{X};\theta)$. Alg.: Iteratively $q^\star \in \text{argmin}_q E(q, \theta_{t-1})$ and $\theta_t \in \text{argmax}_\theta M(q^\star, \theta)$. These can be done in closed form for GMM.

### Density estimation

MLE properties: (1) *Consistency*: $\lim_{n\to\infty} \hat{\theta}_n^{\text{MLE}} = \theta$; (2) *Equivariance*: If $\hat{\theta}$ is the MLE of $\theta$, then $g(\hat{\theta})$ is the MLE of $g(\theta)$; (3) *Asymptotically normal*: In the limit of $n$, $\hat{\theta} - \theta/\sqrt{n}$ converges to $\mathcal{N}(0, \mathcal{I}(\theta)^{-1})$; (4) *Asymptotically efficient*: In the limit of $n$, MLE has smallest variance among unbiased estimators.

Rao-Cramér bound: For any unbiased estimator:

$$\mathbb{V}[\hat{\theta}(y)] \geq \frac{(\frac{\partial}{\partial \theta}b_{\hat{\theta}} + 1)^2}{\mathcal{I}_n(\theta)} + b_{\hat{\theta}}^2,$$

where $\mathcal{I}_n(\theta) \doteq \mathbb{E}_{y|\theta}[(\frac{\partial}{\partial \theta}\log p(y \mid \theta))^2]$ and

---

$b_{\hat{\theta}} \doteq \mathbb{E}_{y|\theta}[\hat{\theta}(y)] - \theta$. If unbiased: $\mathbb{V}[\hat{\theta}(y)] \geq 1/\mathcal{I}_n(\theta)$.
And MLE: $\lim_{n\to\infty}\mathbb{V}[\hat{\theta}^{\text{MLE}}(y)] = 1/\mathcal{I}_n(\theta)$.

### Regression

Minimize loss: $\ell(f) = \frac{1}{n}\sum_{i=1}^n (f(x_i) - y_i)^2$.

**Linear regression:** Assume $Y \mid X = x \sim \mathcal{N}(\beta_\star^\top x, \sigma^2)$. We parameterize $f(x;\beta) = \beta^\top x$.

OLSE: $\hat{\beta} = (\mathbf{X}^\top \mathbf{X})^{-1}\mathbf{X}^\top y$ s.t. $\mathbf{X} \in \mathbb{R}^{n \times d}, y \in \mathbb{R}^n$.
Potential problems:

1. Remove outliers, because linear models are heavily influenced by them;
2. Standardize data, because features on different scales result in unstable matrix inversion;
3. "Curse of dimensionality": In high dimensionality, logistic regression outputs overconfident outputs, due to overestimation of weights;
4. Collinear features result in unstable matrix inversion due to small eigenvalues.

*Risk decomposition*:
$\mathbb{E}[(\hat{f}(X) - Y)^2] = (\mathbb{E}[\hat{f}(X)] - \mathbb{E}[y])^2 + \mathbb{V}[\hat{f}(X)] + \mathbb{V}[y]$.
*Proof*: Use $Y = f(X) + \epsilon$ and show $\mathbb{E}[\epsilon^2] = \mathbb{V}[Y]$ where $\mathbb{E}[\epsilon] = 0$. Then $\pm\mathbb{E}[\hat{f}(X)]$ and finalize.

*Gauss-Markov*: $\mathbb{V}[a^\top \hat{\beta}] \leq \mathbb{V}[a^\top \tilde{\beta}]$ for any $a \in \mathbb{R}^d$ and $\tilde{\beta} = \mathbf{C}y$ for $\mathbf{C} \in \mathbb{R}^{d \times n}$. (OLSE $\hat{\beta}$ is unique min.-var. unbiased linear estimator.) This does not mean it is best, because adding some bias may decrease variance considerably.

**Regularization:** Ridge: Gaussian prior $\beta \sim \mathcal{N}(0, \lambda \mathbf{I})$. LASSO: Laplacian Gaussian $\beta \sim \text{Lap}(0, \lambda \mathbf{I})$. $\ell_1$ results in sparse weights (better interpretation) and the sign of features remain.

**Polynomial regression:** Feature map with all polynomials $\phi(x)$ and perform lin. reg. in this space: $\psi(x;\beta) = \beta^\top \phi(x)$. Problem: Infinitely dimensional $\Rightarrow$ Ill-defined inner product. Solution: Fix by data-dependent scalar, specifically: $\phi(x) = \exp(-\|x\|^2/2)\left[\prod_{i=1}^d x_i^{\alpha_i}/\sqrt{\prod_{i=1}^d \alpha_i!}\right]_{\alpha \in \mathbb{N}^d}$. $\Rightarrow$ RBF kern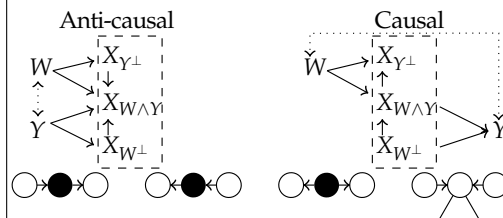el: $\langle\phi(x), \phi(x')\rangle = \exp(-\|x-x'\|^2/2)$. Now compute OLSE in this space: $\hat{\beta} = (\mathbf{\Phi}^\top \mathbf{\Phi})^{-1}\mathbf{\Phi}^\top y$, $\mathbf{\Phi} \in \mathbb{R}^{n \times \infty}$. Problem: Cannot compute $\mathbf{\Phi}^\top \mathbf{\Phi} \in \mathbb{R}^{\infty \times \infty}$. Solution: Rewrite OLSE: $\hat{\beta} = \mathbf{\Phi}^\top(\mathbf{\Phi}\mathbf{\Phi}^\top)^{-1}y$. Prediction only contains kernel evaluations: $\psi(x) = k(x)^\top \mathbf{K}^{-1}y$. Problem: $\mathcal{O}(n^3)$ runtime.

### Causality

Causal fallacies where one might conclude $X$ causes $Y$: (1) Reverse causality: $Y$ causes $X$; (2) Third-cause fallacy: $Z$ causes $X$ and $Y$; (3) Bidirectional causation: $X$ causes $Y$ and $Y$ causes $X$.
*Domain shift*: Test samples are drawn from different distribution than training set.
*Shortcut learning*: Spurious correlation between causal and non-causal features in the training depend on environment.



Necessary conditions for counterfactual invariance: Anti-causal: $f(\mathbf{X}) \perp \mathbf{W} \mid Y$. Causal without selection (but possibly confounded): $f(\mathbf{X}) \perp \mathbf{W}$. Causal without confounded (but possibly selection): $f(\mathbf{X}) \perp \mathbf{W} \mid Y$ as long as $\mathbf{X} \perp Y \mid \mathbf{X}_{W^\perp}, \mathbf{W}$.

### Gaussian processes

Outputs are modeled as $y = \mathbf{X}\beta + \epsilon, \epsilon \sim \mathcal{N}(0, \mathbf{I})$. Thus: $y \mid \mathbf{X}, \beta \sim \mathcal{N}(\mathbf{X}\beta, \sigma^2 \mathbf{I})$. BLR extends linear reg. with prior on $\beta$: $\beta \sim \mathcal{N}(0, \mathbf{\Lambda}^{-1})$.

---

Posterior: $\beta \mid \mathbf{X}, y \sim \mathcal{N}(\tilde{\mu}, \tilde{\mathbf{\Sigma}})$, where

$$\tilde{\mu} = -\frac{1}{\sigma^2}\tilde{\mathbf{\Sigma}}\mathbf{X}^\top y, \quad \tilde{\mathbf{\Sigma}} = \sigma^2\left(\mathbf{X}^\top\mathbf{X} + \sigma^2\mathbf{\Lambda}\right)^{-1}.$$

Joint distribution over outputs (using prior):
$y \mid \mathbf{X} \sim \mathcal{N}\left(0, \mathbf{X}\mathbf{\Lambda}^{-1}\mathbf{X}^\top + \sigma^2\mathbf{I}\right)$. Prediction: $y^\star \mid x^\star, \mathbf{X}, y \sim \mathcal{N}(\mu^\star, \Sigma^\star)$, where $\mu^\star = k^\top(\mathbf{K} + \sigma^2\mathbf{I})^{-1}y$ and $\Sigma^\star = k - k^\top(\mathbf{K} + \sigma^2\mathbf{I})^{-1}k$. Problem: $\mathcal{O}(n^3)$ runtime.

**Kernels:** Kernel $k$ must satisfy symmetry and PSD:

$$\int\int f(x)k(x,x')f(x')dxdx' \geq 0, \quad \forall f \in L_2.$$

Or there exists $\phi$ s.t. $k(x,x') = \phi(x)^\top\phi(x')$.

*Linear kernel*: $k(x,x') = x^\top x'$; *Polynomial kernel*: $k(x,x') = (x^\top x' + 1)^p$; *RBF kernel*: $k(x,x') = \exp(-\|x-x'\|^2/\ell^2)$; *Sigmoid kernel*: $\tanh(\kappa x^\top x') - b$. If $k_1$ and $k_2$ are valid kernels and $c > 0$, then the following are: $k_1 + k_2$, $k_1 \cdot k_2$, $c \cdot k_1$, and $\exp \circ k_1$.

### Uncertainty quantification

**Statistical model validation:** Methods to evaluate $f$ (or algorithm $\mathcal{A}$) that is trained on data $\mathcal{Z}$: *Cross-validation*: Partition $\mathcal{Z} = \bigcup_{k=1}^K \mathcal{Z}_k$ and produce $K$ estimators $\hat{f}^{-k}$ from $\mathcal{Z} \backslash \mathcal{Z}_k$. Then estimate risk by $\mathcal{R}^{\text{CV}}(\mathcal{A}) = \frac{1}{n}\sum_{i=1}^n \ell(y_i, \hat{f}^{-k(i)}(x_i))$, where $k$ maps $i$ to the partition such that $x_i \in \mathcal{Z}_{k(i)}$.

*Bootstrap*: Used for measuring dist. over stat. params. Draw $B$ bootstrap samples $\Rightarrow$ Compute parameter for each $\Rightarrow$ Compute statistics. Can also use for empirical risk: $\hat{\mathcal{R}}^{\text{BS}}(\mathcal{A}) \doteq \frac{1}{n \cdot B}\sum_{b=1}^B \sum_{i=1}^n \ell(y_i, \hat{f}^{*b}(x_i))$. Problem: Overly optimistic. Solution: $\mathcal{R}^{\text{BS}}(\mathcal{A}) \doteq \frac{1}{n}\sum_{i=1}^n \frac{1}{|\mathcal{C}^{-i}|}\sum_{b \in \mathcal{C}^{-i}} \ell(y_i, \hat{f}^{*b}(x_i))$.

Correct for optimism of $\hat{\mathcal{R}}^{\text{BS}}$ by combining with $\mathcal{R}^{\text{BS}}$: $\mathcal{R}^{(0.632)} = 0.368\hat{\mathcal{R}}^{\text{BS}} + 0.632\mathcal{R}^{\text{BS}}$. 0.632 is the prob. that a sample appears at least once in a bootstrap sample of size $n$.

**Uncertainty in linear models:** OLSE has distribution over estimators: $\hat{\beta} \sim \mathcal{N}(\beta^\star, \sigma^2(\mathbf{X}^\top\mathbf{X})^{-1})$. Unbiased estimator of $\sigma^2$: $\hat{\sigma}^2 = \frac{1}{n-d}\sum_{i=1}^n(\hat{\beta}^\top x_i - y_i)$. Then we have $1 - \alpha$ confidence interval for $\beta_j^\star$: $\hat{\beta}_j \pm z_{\alpha/2}\epsilon(\hat{\beta}_j)$, where $z_{\alpha/2} = \Phi^{-1}(\alpha/2)$, $\Phi$ is standard Gaussian CDF, and $\epsilon(\hat{\beta}_j) = \hat{\sigma}^2(\mathbf{X}^\top\mathbf{X})_{jj}^{-1}$.

**Statistical testing:** Null hypothesis: $H_0 : \theta^\star \in \Theta$. Alternative hypothesis $H_1 : \theta^\star \in \Theta$. We are given $n$ samples $x_1, \ldots, x_n \sim p(\cdot \mid \theta^\star)$ and a test statistic $t : \mathcal{X}^n \to \mathbb{R}$. The goal is to find a critical value $c \in \mathbb{R}$ such that $\mathbb{P}(t(X_1, \ldots, X_n) \geq c \mid \theta)$ is low when $\theta \in \Theta_0$ and high when $\theta \in \Theta_1$.

We want to minimize the prob. of choosing $H_1$ when $H_0$ holds (worst possible situation). We quantify this notion of risk as $\alpha_c \doteq \sup_{\theta \in \Theta_0} \mathbb{P}(t(x_1, \ldots, x_n) \geq c \mid \theta)$. Problem: $\alpha_c \to 0$ as $c \to \infty$, so $c^\star \to \infty$ minimizes the risk, but then we never accept $H_1$. Solution: Run test on realization $t(x_1, \ldots, x_n)$ and compute risk of least risky critical value that would incorrectly reject $H_0$: $p = \inf_{c \in \mathbb{R}}\{\alpha_c \mid t(x_1, \ldots, x_n) \geq c\}$. This is the $p$-value:
$p \doteq \sup_{\theta \in \Theta_0} \mathbb{P}(t(X_1, \ldots, X_n) \geq t(x_1, \ldots, x_n) \mid \theta)$.
Intuition: Inverse prob. of $x_{1:n}$ being an outlier.

*Wald*: $W = (\hat{\theta} - \theta_0)^2/\hat{\sigma}^2$; $H_0 : \theta = \theta_0$, $H_1 : \theta \neq \theta_0$.

**Bayesian neural networks:** (S)GD only yields single point estimate of weights $\Rightarrow$ Define prior $\theta \sim \mathcal{N}(0, \sigma^2\mathbf{I})$ and likelihood $p(\mathcal{Z} \mid \theta) = \prod_{x,y \in \mathcal{Z}} p(y \mid x, \theta) \Rightarrow$ Posterior with Bayes rule. Problem: $p(\mathcal{Z})$ is intractable. Solution: Variational inference with isotropic Gaussians and find

$$q^\star \in \underset{\mu,\sigma>0}{\text{argmin}}\, D_{\text{KL}}(\mathcal{N}(\mu, \sigma^2\mathbf{I}) \parallel p(\theta \mid \mathcal{Z}))$$
$$= \underset{\mu,\sigma>0}{\text{argmin}}\, \mathbb{E}_{\theta \sim \mathcal{N}}[F(\mu, \sigma, \theta)],$$

where $F(\mu, \sigma, \theta) = \log\mathcal{N}(\theta; \mu, \sigma^2\mathbf{I}) - \log p(\mathcal{Z} \mid$

$\theta) - \log p(\theta)$. Then, we can apply SGD with the following gradients:
$$\nabla_\mu = \mathbb{E}_\epsilon[\nabla_\theta F(\mu, \sigma, \theta) + \nabla_\mu F(\mu, \sigma, \theta)]$$
$$\nabla_\sigma = \mathbb{E}_\epsilon[\epsilon^\top \nabla_\theta F(\mu, \sigma, \theta) + \nabla_\sigma F(\mu, \sigma, \theta)],$$
where $\epsilon \sim \mathcal{N}(0, I)$ and $\theta = \mu + \sigma\epsilon$.

**Information-based transductive learning:** We are given domain $\mathcal{X}$ that contains safe area $\mathcal{S} \subseteq \mathcal{X}$ and area of interest $\mathcal{A} \subseteq \mathcal{X}$. We have an unknown $f^\star$ that we want to explore within $\mathcal{A}$, but we can only query (noisy) observations in $\mathcal{S}$:
$$y_x = f^\star(x) + \epsilon_x, \quad \mathbb{E}_{\epsilon_x} = 0.$$
We are given a history of points $\mathcal{D}_{n-1}$ and need to compute which point will give the most additional information. ITL selects the next point as:
$$x_n \in \operatorname*{argmax}_{x \in \mathcal{S}} \mathrm{I}(f_\mathcal{A}; y_x \mid \mathcal{D}_{n-1}).$$
If $f \sim GP(\mu, k)$, then
$$\mathrm{I}(f_\mathcal{A}; y_x \mid \mathcal{D}_{n-1}) = \frac{1}{2}\log\frac{\mathbb{V}[y_x \mid \mathcal{D}_{n-1}]}{\mathbb{V}[y_x \mid f_\mathcal{A}, \mathcal{D}_{n-1}]}.$$
*Proof*: Use entropy of Gaussian.

## Convex optimization and SVMs

$$\min \quad f(x), \quad \text{s.t.} \quad g_i(x) = 0, h_j(x) \le 0.$$
where $f$ and $h_j$ are convex and $g_i$ are affine.

Lagrangian: $\mathcal{L}(x, \lambda, \alpha) \doteq f(x) + \sum_{i=1}^n \lambda_i g_i(x) + \sum_{j=1}^m \alpha_j h_j(x)$. Lagrange dual function:
$\theta(\lambda, \alpha) \doteq \inf_{x \in \mathcal{X}} \mathcal{L}(x, \lambda, \alpha)$.

Weak duality: Let $x \in \mathcal{C}$, $\alpha \ge 0$, then $\theta(\lambda, \alpha) \le f(x)$. Thus: $\max_{\lambda, \alpha \ge 0} \theta(\lambda, \alpha) \le \min_{x \in \mathcal{C}} f(x)$.

If there is a Slater point (exists $x \in \mathcal{C}$ such that $h_j(x) < 0$ for all $j$) then strong duality: $\max_{\lambda, \alpha \ge 0} \theta(\lambda, \alpha) = \min_{x \in \mathcal{C}} f(x)$.

If all $g_i$ and $h_j$ are differentiable, KKT conditions provide necessary (and sufficient for convex programming) conditions for strong duality:
$$\alpha_j^\star h_j(x^\star) = 0, \quad \nabla_x \mathcal{L}(x^\star, \lambda^\star, \alpha^\star) = 0.$$
Or, condition 2: $x^\star \in \operatorname{argmin}_{x \in \mathcal{X}} \mathcal{L}(x, \lambda^\star, \alpha^\star)$.

**Support vector machine:** We want to linearly separate a dataset with maximum margin $\Rightarrow$ Model as convex program with constraint for each data point: $f[w, b](x, y) = y(w^\top x + b) \ge \epsilon > 0$.

Margin ($x^+$ and $x^-$ are support vectors):
$$2 \cdot m(w, b) = \|\operatorname{proj}_w(x^+) - \operatorname{proj}_w(x^-)\|$$
$$= |\bar{w}^\top(x^+ - x^-)|.$$
Ill-posed problem because infinite number of solutions $\Rightarrow$ Only one solution satisfies
$$w^\top x^+ + b = 1, \quad w^\top x^- + b = -1.$$
Then, $m(w, b) = 1/\|w\|$:
$$\min \quad \frac{1}{2}\|w\|^2, \quad \text{s.t.} \quad 1 - y_i(w^\top x_i + b) \le 0.$$
$w^\star = \sum_{i=1}^n \alpha_i^\star y_i x_i$, $b^\star = -\frac{1}{2}(w_\star^\top x^+ + w_\star^\top x^-)$, where $\alpha^\star$ is the dual solution.

**SVM variations:** Soft-margin introduces slackness in case data is not linearly separable:
$$\text{minimize} \quad \frac{1}{2}\|w\|^2 + C\sum_{i=1}^n \xi_i$$
$$\text{subject to} \quad y_i(w^\top x_i + b) \ge 1 - \xi_i, \quad \xi_i \ge 0.$$
Optimal slackness parameters:
$$\xi_i^\star = \max\{0, 1 - y_i(w_\star^\top x_i + b^\star)\}.$$
If data is not linearly separable, use features and their kernels: $w_\star^\top \phi(x) = \sum_{i=1}^n \alpha_i^\star y_i k(x_i, x)$.

We can generalize the margin notion to multiclass by introducing weights $w_z$ per class. The margin is defined as the maximum $m \in \mathbb{R}$ s.t.
$$m \le (w_{z_i}^\top y_i + b_{z_i}) - \max_{z \ne z_i}\{w_z^\top y_i + b_z\}.$$
New optimization problem:
$$\min \quad \frac{1}{2}\|w\|^2 = \frac{1}{2}\sum_{z=1}^k \|w_z\|^2$$
$$\text{s.t.} \quad (w_{z_i}^\top y_i + b_{z_i}) - \max_{z \ne z_i}\{w_z^\top y_i + b_z\} \ge 1.$$

---

Structural SVMs can have infinitely many classes. So, we need to define a joint feature map $\psi$ such that $f_w(x, y) = w^\top \psi(x, y)$. This is used to perform classification: $c(x) = \operatorname{argmax}_{y \in \mathcal{Y}} f_w(x, y)$. We need to construct an algorithm to efficiently compute this argmax and an algorithm to compute the max in the below optimization problem. Some structures are closer than others $\Rightarrow$ Introduce a loss function $\Delta$:
$$\min \quad \frac{1}{2}\|w\|^2 \quad \text{s.t.} \quad w^\top \psi(x_i, y_i)$$
$$- \max_{y \ne y_i}\{w^\top \psi(x_i, y) + \Delta(y, y_i)\} \ge 0.$$

## Ensembles

Average $B$ estimators into $\hat{f} \Rightarrow$ avg. bias and:
$$\mathbb{V}[\hat{f}] = \frac{1}{B^2}\sum_{b=1}^B \mathbb{V}[\hat{f}_b] + \frac{1}{B^2}\sum_{b=1}^B \sum_{b' \ne b}^B \operatorname{Cov}(\hat{f}_b, \hat{f}_{b'}).$$
If the covariances are low, the variance is significantly decreased while the bias remains the same.

**Bagging:** $B$ times take a bootstrap sample and train a classifier. This works well because covariances are small due to using different subsets for training and the variances are similar because each subsample behaves similarly on average.

Random forests do this with (very deep) decision trees. Very deep because they have low bias and high variance, which is reduced by ensembling.

**AdaBoost:** AdaBoost reduces cov. by using a different weighting for each estimator. The weights are determined by error of prev. classifiers.
$$w_i^{(b+1)} = w_i^{(b)}\exp(\alpha_b \mathbb{1}\{c_b(x_i) \ne y_i\})$$
$$\alpha_b = \log(1 - \epsilon_b/\epsilon_b)$$
$$\epsilon_b = \sum_{i=1}^n \frac{w_i^{(b)}}{\sum_{j=1}^n w_j^{(b)}}\mathbb{1}\{c_b(x_i) \ne y_i\}.$$
Final classifier: $\hat{c}(x) = \operatorname{sgn}(\sum_{b=1}^B \alpha_b c_b(x))$.
It can be shown that AdaBoost fits an additive model in base learners optimizing the exponential loss $\mathbb{E}[\exp(-yf(x))]$ via Newton-like updates.

## Stable diffusion

**Diffusion models:** Iteratively denoise. Continuous:
$$dx_t^+ = \mu(x_t, t)dt + \sigma(x_t, t)d\omega_t$$
$$dx_t^- = [\mu(x_t, t) - \sigma^2(x_t, t)\nabla_x \log p_t(x_t)]dt + \sigma(x_t, t)d\bar{\omega}_t.$$
DDPM scheduler: $x_{t+1} = \sqrt{1 - \beta_t}x_t + \sqrt{\beta_t}\epsilon$, where $\epsilon \sim \mathcal{N}(0, I)$. Backward process:
$x_{t-1} = \frac{1}{\sqrt{\alpha_t}}(x_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}}\epsilon_\theta(x_t, t)) + \sqrt{\beta_t}z$,
where $\alpha_t = 1 - \beta_t$ and $\bar{\alpha}_t = \prod_{\tau=1}^t \alpha_t$.
Diffusion models are trained by sampling $x_0 \sim p_0$, $t \sim \operatorname{Unif}(\{1, \ldots, T\})$, $\epsilon \sim \mathcal{N}(0, I)$ and performing gradient step on $\ell = \|\epsilon - \epsilon_\theta(x_t, t)\|^2$.

## Non-parametric Bayesian methods

Beta distribution ($x \in [0, 1], \alpha, \beta > 0$): $\operatorname{Beta}(x; \alpha, \beta) \propto x^{\alpha-1}(1 - x)^{\beta-1}$. Dirichlet generalizes ($x \in \Delta^{n-1}$): $\operatorname{Dir}(x; \alpha) \propto \prod_{k=1}^n x_k^{\alpha_k - 1}$.
Problem: Need to know $K$ (#clusters) beforehand.
A Dirichlet process $DP(\alpha, H)$ is a distribution over probability distributions on a space $\Theta$, where $\alpha$ is a concentration parameter. A sample $G \sim DP(\alpha, H)$ is a function $G : \Theta \to \mathbb{R}_{\ge 0}$ such that $\int_\Theta G(\theta)d\theta = 1$. For every partition $(T_1, \ldots, T_k)$ of $\Theta$ and $G \sim DP(\alpha, H)$, we have
$$(G(T_1), ., G(T_k)) \sim \operatorname{Dir}(\alpha H(T_1), ., \alpha H(T_k)).$$
Dir can be sampled recursively by stick-breaking:
$$\beta_i \sim \operatorname{beta}\left(\alpha_i, \prod_{k=i+1}^K \alpha_k\right), \quad \rho_i = \beta_i \prod_{j=1}^{i-1}(1 - \beta_j)$$
$$(\rho_{i+1}, \ldots, \rho_K) \sim \operatorname{Dir}(\alpha_{i+1}, \ldots, \alpha_K).$$

---

Still limited to fixed $K$. GEM distribution fixes this by fixing $\alpha$ such that $\beta_i \sim \operatorname{Beta}(1, \alpha)$ for all $i$. Recursion: $\beta_i \sim \operatorname{Beta}(1, \alpha)$ and
$$\rho_i = \beta_i \prod_{j=1}^{i-1}(1 - \beta_j), \quad \rho_K = \beta_k\left(1 - \sum_{i=1}^{K-1}\rho_i\right).$$
Keep sampling cluster probs until satisfied.
If $(\rho_1, \rho_2, \ldots) \sim \operatorname{GEM}(\alpha)$ and $\theta_k \sim H$, then this is sample from $DP(\alpha, H)$: $G(\theta) = \sum_{k=1}^\infty \rho_k \delta_{\theta_k}(\theta)$.

**Chinese restaurant process:**
$$P(n + 1 \text{ joins table } \theta \mid \mathcal{P}) = \begin{cases} \frac{|\theta|}{\alpha + n} & \theta \in P \\ \frac{\alpha}{\alpha + n} & \text{else.} \end{cases}$$
Probability of partition $\mathcal{P}$ can be written as
$$P(\mathcal{P}) = \alpha^{|\mathcal{P}|}\frac{\alpha!}{(N + \alpha)!}\prod_{\tau \in \mathcal{P}}(|\tau| - 1)!.$$
Problem is exchangeable. $\mathbb{E}[|\mathcal{P}|] \in \mathcal{O}(\alpha \log N)$.

**DPMM:** Assume $\Theta = \mathbb{R}$ with $\mu \in \Theta$ corresponding to $\mathcal{N}(\mu, \sigma)$ for fixed $\sigma > 0$ and $H = \mathcal{N}(\mu_0, \sigma_0)$ for fixed $\mu_0, \sigma_0$. DPMM: Cluster probs are sampled from GEM: $(\rho_1, \rho_2, \ldots) \sim \operatorname{GEM}(\alpha)$. Cluster centers are sampled from base measure: $\mu_1, \mu_2, \ldots \sim \mathcal{N}(\mu_0, \sigma_0)$. Clusters are assigned: $z_i \sim \operatorname{Cat}(\rho_1, \rho_2, \ldots), \forall i \in [n]$. Data points are sampled: $x_i \sim \mathcal{N}(\mu_{z_i}, \sigma), \forall i \in [n]$. This process is exchangeable. To fit a DPMM, we use a collapsed Gibbs sampling formulation:
$p(z_i = k \mid z^{-i}, x, \alpha, \mu) \propto p(z_i = k \mid z^{-i}, \alpha)p(x_i \mid x^{-i}, z_i = k, z^{-i}, \mu)$. Prior is as in CRP:
$$p(z_i = k \mid z^{-i}, \alpha) = \begin{cases} \frac{N_k^{-i}}{\alpha + N - 1} & \text{existing } k \\ \frac{\alpha}{\alpha + N - 1} & \text{else.} \end{cases}$$
Likelihood (right term) is cond. on cluster $k$:
$$\ell = \begin{cases} p(x_i \mid x_k^{-i}, \mu) = \frac{p(x_i x_k^{-i} \mid \mu)}{p(x_k^{-i} \mid \mu)} & \text{existing } k \\ p(x_i \mid \mu) & \text{else.} \end{cases}$$

## Statistical learning theory

**PAC learning:** *Definition*: A learning algorithm $\mathcal{A}$ can learn a concept $c \in \mathcal{C}$ if there exists $\operatorname{poly}(\cdot, \cdot, \cdot)$ such that for any distribution $p$ on $\mathcal{X}$ and $\epsilon, \delta \in (0, 1/2)$, if $\mathcal{A}$ receives a sample of size $n \ge \operatorname{poly}(1/\epsilon, 1/\delta, \operatorname{size}(c))$, then $\mathcal{A}$ outputs $\hat{c}$ such that $\mathbb{P}(\mathcal{R}(\hat{c}) \le \epsilon) \ge 1 - \delta$. This probability is taken over the randomness of $\mathcal{Z}$ and $\mathcal{A}$.
$\mathcal{C}$ is PAC learnable from $\mathcal{H}$ if there is an $\mathcal{A}$ that can learn any $c \in \mathcal{C}$.
If $\mathcal{A}$ runs polynomial in only $1/\delta$ and $1/\epsilon$, then $\mathcal{C}$ is efficiently PAC learnable.
In the stochastic setting, $y$ is also random and not deterministically decided by a concept $c \in \mathcal{C}$. Now the criterium is $\mathbb{P}_{\mathcal{Z} \sim p}(\mathcal{R}(\hat{c}) - \inf_{c \in \mathcal{C}} \mathcal{R}(c) \le \epsilon) \ge 1 - \delta$.

**Vapnik-Chervonenkis:** VC dimension is the cardinality of the largest set of points that $\mathcal{C}$ can shatter.
*Vapnik and Chervonenkis*: Assume a finite concept class and $\mathcal{R}(c^\star) = 0$ and define $c_n^\star \in \{c \in \mathcal{C} \mid \hat{\mathcal{R}}_n(c) = 0\}$. Then, for every $n \in \mathbb{N}$ and $\epsilon > 0$,
$$\mathbb{P}(\mathcal{R}(\hat{c}_n^\star) > \epsilon) \le |\mathcal{C}|\exp(-n\epsilon).$$
And: $\mathbb{E}[\mathcal{R}(\hat{c}_n^\star)] \le \frac{1 + \log|\mathcal{C}|}{n}$.
*VC inequality*: $\mathbb{P}(\mathcal{R}(\hat{c}_n^\star) - \inf_{c \in \mathcal{C}} \mathcal{R}(c) > \epsilon) \le \mathbb{P}(\sup_{c \in \mathcal{C}} |\hat{\mathcal{R}}_n(c) - \mathcal{R}(c)| > \frac{\epsilon}{2})$.
*Hoeffding*: Let $X_i \in [a_i, b_i]$ be i.i.d. and $S_n = \sum_{i=1}^n X_i$, then for any $t > 0$,
$$\mathbb{P}(S_n - \mathbb{E}[S_n] \ge t) \le \exp\left(-\frac{2t^2}{\sum_{i=1}^n (b_i - a_i)^2}\right).$$
Same bound for $\le -t$.
As a result: $\mathbb{P}(\tilde{S}_n - \mathbb{E}[\tilde{S}_n] \ge \epsilon) \le \exp\left(-\frac{2n\epsilon^2}{\sum_{i=1}^n (b_i - a_i)^2/n}\right)$, where $\tilde{S}_n = S_n/n$.
Assume $|\mathcal{C}| \le N$, then for all $\epsilon > 0$,
$$\mathbb{P}(\sup_{c \in \mathcal{C}} |\hat{\mathcal{R}}_n(c) - \mathcal{R}(c)| > \epsilon) \le 2N\exp(-2n\epsilon^2).$$
We can deal with infinite $|\mathcal{C}|$ by representing hypotheses by the classifications that they yield. Or measuring the VC dimension of $\mathcal{C}$.