# Left Column

**Regression:** $N$: data-size, $D$: dim (Features), $D+1$: Model param. with offset.
- relates $x_n$ to $y_n$. Regression Func°: $y_n \approx f(x_n) = f_w(x_n)$ → 2 goals: Prediction, Interpr°

**Linear Regression:** Univariate: $y_n \approx f(x_n) = w_0 + w_1 x_{n1}$

Multivariate: $y_n \approx f(x_n) = w_0 + w_1 x_{n1} + \dots + w_D x_{nD} = w_0 + x_n^\top \binom{w_1}{w_D} = \tilde{x}_n^\top \tilde{w}$

Note: If $D > N$. We say the model is overparametrized → Solution: Regularization $= \binom{1}{x_n} \tilde{w} \in \mathbb{R}^{D+1}$

**Cost Func°:** $\mathcal{L}(w) = \frac{1}{N} \sum_{n=1}^{N} \text{Loss}(e_n)$, $e_n = y_n - f_w(x_n)$ | 2 desirables pr°prties: (sum over training samples)
1) Symmetric around 0  2) Penalize large and huge mistakes the same (outliers robust)

$MSE(w) := \frac{1}{N} \sum_{n=1}^{N} [y_n - f_w(x_n)]^2$    $MAE(w) := \frac{1}{N} \sum_{n=1}^{N} |y_n - f_w(x_n)|$ (more robust to outliers than MSE)

**Convexity:** A function $h(u)$ with $u \in \mathbb{R}^D$ is convex if for any $u,v \in \mathbb{R}^D$ and any $0 \leq \lambda \leq 1$: $h(\lambda u + (1-\lambda)v) \leq \lambda h(u) + (1-\lambda)h(v)$ (desirable computational property)

- A strict. convex function has a unique global min. For convex functions, every local min is a global min
- Sum of Convex functions are convex    $f: \mathbb{R}^D \to \mathbb{R}$ convex iff $\nabla^2 f = \text{PSD}$
- Local opt: $\mathcal{L}(w^*) \leq \mathcal{L}(w)$, $\forall w$ with $\|w - w^*\| < \epsilon$ (strict) | Global: $\mathcal{L}(w^*) \leq \mathcal{L}(w)$, $\forall w \in \mathbb{R}^D$

**LS:** Solving Lin Reg with MSE analytically. 2 steps for global opt: 1) Show convexity 2) $\nabla \mathcal{L}(w) = 0$

Normal eq: $\mathcal{L}(w) = \frac{1}{2N} \sum_{n=1}^{N} (y - x_n^\top w)^2 = \frac{1}{2N} e^\top e$ with $e = Y - Xw$
1) Cost fct is convex bc it's combination of convex fcts | Cost: $\Theta(ND^2 + D^3)$

2) $\nabla \mathcal{L}(w) = -\frac{1}{N} X^\top e = 0 \Rightarrow X^\top X w = X^\top Y \Rightarrow w^* = (X^\top X)^{-1} X^\top Y$ iff Gram matrix $X^\top X \in \mathbb{R}^{D \times D}$ is invertible (if rank(X) = D) Prediction on unseen data:
→ $\hat{y}_m = x_m^\top w^* = x_m^\top (X^\top X)^{-1} X^\top Y$. • if $D > N$: rank(X) < D • if $D \leq N$, but some col X:d are (nearly) collinear: ill-cond. ⇒ Solution: Linear Solver: $X^\top X w = X^\top Y$

**Underfitting:** Can't find fit, too limited   **Overfitting:** Fits noise too, too rich

Augment linear models: $y_n \approx w_0 + w_1 x_n + w_2 x_n^2 + \dots + w_M x_n^M = \phi(x_n)^\top w$

Avoid overfit: Increase $N$, keep $M$ fixed (add data, same model complexity)

**MLE:** Gaussian RV in $\mathbb{R}^D$: $\mathcal{N}(y \mid \mu, \Sigma) = (2\pi)^{-D/2} \det(\Sigma)^{-1/2} \exp[-0.5(y-\mu)^\top \Sigma^{-1}(y-\mu)]$ ($\Sigma$ psd)
in $\mathbb{R}$: $P(y \mid \mu, \sigma^2) = \mathcal{N}(y \mid \mu, \sigma^2) = (2\pi \sigma^2)^{-1/2} \exp[-(y-\mu)^2/2\sigma^2]$

Note: When 2 RV indep $P(x,y) = P(x) P(y)$

**Probabilistic model for LS:** assume data $y_n = x_n^\top w + \epsilon_n$ with $\epsilon_n \sim \mathcal{N}(0, \sigma^2)$ (unknown)

Likelihood: $p(y \mid X, w) = \prod_{n=1}^{N} p(y_n \mid x_n, w) = \prod_{n=1}^{N} \mathcal{N}(y_n \mid x_n^\top w, \sigma^2)$ (Independence, assumption on noise) Goal: to max it, and it's equivalent as min $\mathcal{L}_{MSE}$

Log $\mathcal{L}_{LL}(w) = \log p(y \mid X, w) = -\frac{1}{2\sigma^2} \sum (y_n - x_n^\top w)^2 + \text{cnst.}$

$\arg\min_w \mathcal{L}_{MSE}(w) = \arg\max_w \mathcal{L}_{LL}(w)$

MLE can also be interpreted as find-ing the model under which the observed data is most likely to have been generated from (probabilistically)

**Properties:**
- MLE is a sample approximation to the expected log-likelihood: $\mathcal{L}_{LL}(w) \approx \mathbb{E}_{p(y,x)}[\log p(y \mid x, w)]$
- MLE is consistent, i.e., it will give us the correct model assuming that we have a sufficient amount of data.
- $w_{MLE} \xrightarrow{P} w_{true}$   MLE is efficient, i.e. it achieves the Cramer-Rao lower bound

**Replace Gaussian with Laplace:** $P(y_n \mid x_n, w) = (2b)^{-1} \exp[\frac{-1}{b} |y_n - x_n^\top w|]$

**Generalization:** Data Model: unknown distrib $\mathcal{D}$ with range $X \times Y$

Dataset: $S = \{(x_n, y_n)\}_{n=1}^{N} \sim \mathcal{D}$ i.i.d. Learning algo: $A(S) = f_s$ (output) (S: input)

Expect. error: $L_\mathcal{D}(f) = \mathbb{E}_{(x,y) \sim \mathcal{D}}[\ell(y, f(x))]$ (loss) But $\mathcal{D}$ is unknown | Empirical error: (LLN): $L_S \xrightarrow{|S| \to \infty} L_\mathcal{D}$

$L_S(f) = \frac{1}{|S|} \sum_{(x,y) \in S} \ell(y_n, f(x_n))$ Training error: $L_S(f) = \frac{1}{|S|} \sum_{(x_n, y_n) \in S} \ell(y_n, f_s(x_n))$
But could diverge from $L_\mathcal{D}$ bc of overfitting

Split data (test error): $L_{test}(f_{S_{train}}) = \frac{1}{|S_{test}|} \sum_{(x,y) \in S_{test}} \ell(y_n, f_{S_{train}}(x_n)) \approx L_\mathcal{D}(f_{S_{train}})$

Generalization error: $P[|L_\mathcal{D}(f) - L_{test}(f)| \geq \sqrt{\frac{(b-a)^2 \ln(2/\delta)}{2|S_{test}|}}] \leq \delta$ (Error decreases as $\sqrt{1/(|S_{test}|^N)}$) (# of test points) (Gener. Gap) ($a,b$) bounds on the loss (Upper bound)

Hoeffding's: $P[|\frac{1}{N} \sum_{n=1}^{N} \Theta_n - E(\Theta)| \geq \varepsilon] \leq 2e^{-2N\varepsilon^2/(b-a)^2}$ 1) With $\Theta = \ell(y_n, f(x_n))$
2) $\delta = 2e^{-2|S_{test}|\varepsilon^2/(b-a)^2} \to \varepsilon = \sqrt{\frac{(b-a)^2 \ln(2/\delta)}{2|S_{test}|}}$

**Model Selection:** Split data, train K times for each value, compute error.
Bound: $P[\max_k |L_\mathcal{D}(f_k) - L_{test}(f_k)| \geq \sqrt{\frac{(b-a)^2 \ln(2K/\delta)}{2|S_{test}|}}] \leq \delta$ (same $\Theta$)
Let $k^* = \arg\min_k L_\mathcal{D}(f_k)$ and $\hat{k} = \arg\min_k L_{test}(f_k)$ $P[L_\mathcal{D}(f_{\hat{k}}) \geq L_\mathcal{D}(f_{k^*}) + 2\sqrt{\frac{(b-a)^2 \ln(2K/\delta)}{2|S_{test}|}}] \leq \delta$

**K-fold Cross-Validation:** 1. Randomly partition data into K groups.
2. Train K times, leaving 1 group for test and K-1 for train.
3. Average K results. → Unbiased estimate of the gen. error and its var.

**Logistic reg:** model probab's $\{0,1\}$ $\sigma(z) = (e^z)(1+e^z)^{-1}$, $1 - \sigma(z) = (1+e^z)^{-1}$ (Robust to outliers and unbalanced data)
$\sigma'(z) = \sigma(z)(1-\sigma(z))$, $P(1|x) = \sigma(x^\top w + b)$, $P(0|x) = 1 - P(1|x)$ if $P(1|x) \geq \frac{1}{2}$ → class 1

**MLE:** assume $x \perp w$: $\mathcal{L}(w) \propto \prod_{n=1}^{N} \sigma(x_n^\top w)^{y_n} (1 - \sigma(x_n^\top w))^{1-y_n}$ | LL: $w_* = \arg\min L(w)$
$:= \frac{1}{N} \sum_{n=1}^{N} -y_n x_n^\top w + \log(1 + e^{x_n^\top w})$ Note: if $\{0,1\}$: $\ell(y, g(x)) = -yg(x) + \log(1 + \exp(g(x)))$ if $\{-1,1\}$:
$\log(1 + \exp(-yg(x)))$ $\Rightarrow \nabla L(w) = \frac{1}{N} \sum_{n=1}^{N} (\sigma(x_n^\top w) - y_n) x_n = \frac{1}{N} X^\top (\sigma(Xw) - Y)$ → No CS
solution but convex! | GD: slow $\Theta(N)$ SGD: Fast but Converges slow

**Newton's:** $L(w) \sim L(w_t) + \nabla L(w_t)^\top (w - w_t) + \frac{1}{2}(w - w_t)^\top \nabla^2 L(w_t)(w - w_t)$
$w_{t+1} = w_t - \gamma_t \nabla^2 L(w_t)^{-1} \nabla L(w_t)$ → intensive! Regularized: issue if data lin separable $\Rightarrow \frac{1}{N} \sum_{n=1}^{N} -y_n x_n^\top w + \log(1 + e^{x_n^\top w}) + \frac{\lambda}{2} \|w\|_2^2$

# Right Column

**Optimization:** $w^* = \arg\min \mathcal{L}(w)$ s.t. $w \in \mathbb{R}^D$ → but very sensitive to ill-conditioning → solution: normalize

Smooth: 1) **Gridsearch:** Get the lowest loss over all W. Brute-force | Exponential complexity | No guarantees to find an optimum

2) **Gradient descent:** A gradient (at a point) is the slope of the tangent to the function (at that point). It points to the direction of largest increase of the function.
Vector: $\nabla \mathcal{L}(w) := [\frac{\partial \mathcal{L}(w)}{\partial w_1}, \dots, \frac{\partial \mathcal{L}(w)}{\partial w_D}]^\top \in \mathbb{R}^D$ | update rule: $w^{(t+1)} := w^{(t)} - \gamma \nabla \mathcal{L}(w^{(t)})$ $\gamma \geq 0$ is the step-size or LR (popular: $\gamma = \frac{cst}{\sqrt{t}} \to$ smaller and smaller...)
$\to \mathcal{L} = \frac{1}{2N} \sum (-2(y_n - w_0)) = w_0 - \bar{y} = 0$
Example: GD for 1-param and MSE: $w_0^{(t+1)} := (1 - \gamma) w_0^{(t)} + \gamma \bar{y}$ (Converges for $\gamma \in (0,2)$)

GD for linear MSE: $Y = \begin{bmatrix} y_1 \\ y_n \end{bmatrix}$, $X = \begin{bmatrix} x_{11} \cdots x_{1N} \\ x_{N1} \cdots x_{ND} \end{bmatrix}$, $e = Y - Xw$ $\mathcal{L}(w) = \frac{1}{2N}(y - x_n^\top w)^2 = \frac{1}{2N} e^\top e$ (For MAE: $\nabla \mathcal{L}(w) = -\frac{1}{N} X^\top \text{sign}(e)$)
→ $\nabla \mathcal{L}(w) = -\frac{1}{N} X^\top e$ Cost of GD per step: $\Theta(ND)$ (Loss of 1 data point)

3) **SGD:** 1) Take 1 random $n$  2) $w^{(t+1)} := w^{(t)} - \gamma \nabla \mathcal{L}_n(w^{(t)})$ (Noisy) $\mathbb{E}_n[\nabla \mathcal{L}_n] = \sum_{n=1}^{N} \frac{1}{N} \nabla \mathcal{L}_n(w)$ Cost of SGD per step $\Theta(D)$
cheap and unbiased estimate of the gradient: $\mathbb{E}[\nabla \mathcal{L}_n(w)] = \nabla \mathcal{L}(w)$ $= \nabla \frac{1}{N} \sum \mathcal{L}_n(w) = \nabla \mathcal{L}(w)$ (stable)

4) **Mini-batch SGD:** $w^{(t+1)} := w^{(t)} - \gamma g$ with $g = \frac{1}{|B|} \sum_{n \in B} \nabla \mathcal{L}_n(w^{(t)})$, where B is a subset. → B=1: Classic SGD
B=N: classic GD → g easily parallelized

4.1) **SGD with momentum:** $w^{(t+1)} := w^{(t)} - \gamma m^{(t+1)}$ with $m^{(t+1)} := \beta_1 m^{(t)} + (1-\beta_1) g$

4.2) **ADAM:** $w_i^{(t+1)} := w_i^{(t)} - \frac{\gamma}{\sqrt{v_i^{(t+1)}}} m_i^{(t+1)} \forall i$ with $m^{(t+1)} := \beta_1 m^{(t)} + (1-\beta_1) g$ $v_i^{(t+1)} := \beta_2 v_i^{(t)} + (1-\beta_2)(g_i)^2 \forall i$
- faster forgetting of older weights
- is a momentum variant of Adagrad
- coordinate-wise adjusted learning rate
- strong performance in practice, e.g. for self-attention networks

4.3) **Sign-SGD:** $w_i^{(t+1)} := w_i^{(t)} - \gamma \text{sign}(g_i) \forall i$ (only use the sign (one bit) of each gradient entry → communication efficient for distributed training (but convergence issue)) Sub

**Non-smooth:** (Vector where function is above is a subgradient) if $\mathcal{L}$ convex and differentiable: $g = \nabla \mathcal{L}(w)$ Cost of GD per step: $\Theta(ND)$

1) **Subgradient $\partial$:** Vector $g \in \mathbb{R}^D$ such that $\mathcal{L}(u) \geq \mathcal{L}(w) + g^\top(u - w) \forall u$
$\Rightarrow w^{(t+1)} := w^{(t)} - \gamma g$ Trick: if $\mathcal{L}(w) = h(q(w))$ with $h$ non-diff. and $q$ diff, $\to$ subg of $\mathcal{L}$ at $w$ is $g \in \partial h(q(w)) \cdot \nabla q(w)$ (set of subgr.) Cost of SGD per step $\Theta(D)$

2) **S(S)GD:** $g = $ subgradient to $\mathcal{L}_n(w)$

**Conclusion:** W local min if $\nabla \mathcal{L}(w) = 0$ (critical point) and $\nabla^2 \mathcal{L}(w) \succ 0$ (positive definite)

**Regularization:** Penalize complex models via: $\min_w \mathcal{L}(w) + \Omega(w)$ where $\Omega(w)$ is the regularizer ($\lambda \to 0$ (no reg) → risk overfitting; $\lambda \to \infty$ (strong reg) → risk underfitting)

**$L_2$:** $\Omega(w) = \lambda \|w\|_2^2 = \lambda \sum_i w_i^2$ **Ridge:** $L_2 + MSE$: $\min_w \frac{1}{2N} \sum_{n=1}^{N}(y_n - x_n^\top w)^2 + \lambda \|w\|_2^2$
Closed-form sol: $w_{ridge}^* = (X^\top X + \lambda' I)^{-1} X^\top y$ → always invertible thanks to $\lambda' = 2N\lambda$ ($\lambda'$ lifts the EV so EV are at least $\lambda'$) Proof: $X^\top X + \lambda' I = U \Sigma U^\top + \lambda' U I U^\top = U(\Sigma + \lambda' I) U^\top$
→ Fight ill-conditioning.

**$L_1$:** $\Omega(w) = \lambda \|w\|_1 = \lambda \sum_i |w_i|$ **Lasso:** $L_1 + MSE$: $\min_w \frac{1}{2N} \sum_{n=1}^{N}(y_n - x_n^\top w)^2 + \lambda \|w\|_1$ More powerful. Makes $w$ sparse → few non-0 components
Note: $\{w : \|y - Xw\|^2 = \alpha\}$ is an ellipsoid. These are likely to touch a corner of the "$L_1$-ball" → sparsity.

**Notes:** gradient of $\|w\|^2 = 2w$   subgradient $\|w\|_1 = (\text{sign}(w_1), \dots, \text{sign}(w_D))$

**Bias-Variance:** $Y = f(x) + \varepsilon$. Expected error: $\mathbb{E}_{(x,y)}[((Y - f_S(x))^2]$ → True at every point: $L(f_S) = \mathbb{E}_{\varepsilon \sim b_\varepsilon}[(f(x_0) + \varepsilon - f_S(x_0))^2]$
$\mathbb{E}_{S \sim \mathcal{D}}[L(f_S)] = \mathbb{E}_{S \sim \mathcal{D}, \varepsilon \sim \mathcal{D}_\varepsilon}[(f(x_0) + \varepsilon - f_S(x_0))^2] = \underbrace{\text{Var}_{\varepsilon \sim \mathcal{D}_\varepsilon}[\varepsilon]}_{\text{Noise var}} + \underbrace{(f(x_0) - \mathbb{E}_{S \sim \mathcal{D}}[f_S(x_0)])^2}_{\text{Bias}} + \underbrace{\mathbb{E}_{S \sim \mathcal{D}}[(f_S(x_0) - \mathbb{E}_{S \sim \mathcal{D}}[f_S(x_0)])^2]}_{\text{Var}}$ (can't go below noise var)

**Classification:** unlike reg. $Y \in$ discrete set: $y \in \{C_1, \dots, C_K\}$ Binary $K=2$ **Classifier** divides space into classes

**KNN:** Pros § No opt:, easy, good for low dim § Cons § slow, bad for high dim, distance? § Max-Margin: HP which max margin **Nonlinear classif:** Feature aug. or Kernels **Binary Classif goal:** $\min L_\mathcal{D}(f) = \mathbb{E}_\mathcal{D}(1_{y \neq f(x)}) = P_\mathcal{D}(Y \neq f(x))$

**Bayes classif:** $f^* = \arg\min L_\mathcal{D}(f) \to f^*(x) = \arg\max_{y \in \{-1,1\}} P(Y = y \mid X = x)$ unattainable bc $\mathcal{D}$ not known

2 Classes of classif algo 1) Non-Param: local Avg (KNN) 2) Param: min ERM Problem: Not convex. → Replace $\mathcal{L}$ and $f$ by convex fct.

**ERM:** min ER instead True Risk: $\min_{f: X \to Y} L_{train}(f) := \frac{1}{N} \sum_{n=1}^{N} 1_{f(x_n) \neq y_n} = \frac{1}{N} \sum_{n=1}^{N} y_n f(x_n) \leq 0$

**Kernel:** Ridge: $\min_w \frac{1}{2N} \sum_{n=1}^{N}(y_n - w^\top x_n)^2 + \frac{\lambda}{2} \|w\|^2 \to$ 2 Solutions $w_* = (\frac{1}{N} X^\top X + \lambda I_d)^{-1} X^\top y$ $w_* = \frac{1}{N} X^\top (\frac{1}{N} XX^\top + \lambda I_N)^{-1} y = X^\top \alpha_*$ | $w_* \in \text{span}(x_1, \dots, x_N)$ (Cov $O(d^3 + Nd^2)$) (Kernel $O(N^3 + dN^2)$)

**Representer theorem:** $w_* \in \arg\min_w \frac{1}{N} \sum_{n=1}^{N} \ell(x_n^\top w, y_n) + \frac{\lambda}{2} \|w\|^2$ for ANY loss and min (better than LS) $\Rightarrow w_* = X^\top \alpha_*$ → solution in

For ridge ($\mathcal{L} = \|y - Xw\|^2$): alt formula $\alpha_* = \arg\min_\alpha \frac{1}{N} \alpha^\top (\frac{1}{N} XX^\top + \lambda I_N) \alpha - \alpha^\top y$ | **Kernel matrix:** $K = XX^\top = \begin{bmatrix} x_1^\top x_1 & x_1^\top x_N \\ x_N^\top x_1 & x_N^\top x_N \end{bmatrix} = (x_i^\top x_j)_{i,j} \in \mathbb{R}^{N \times N}$

With feature map $\mathbb{R}^d \to \mathbb{R}^D$: $K = \Phi \Phi^\top = (\Phi(x_i)^\top \Phi(x_j))_{i,j} \in \mathbb{R}^{N \times N}$ Problem: if $d \ll D$, $\phi(x)^\top \phi(x)$ cost $\Theta(D)$ → Expensive!

**Kernel Trick:** K Fct $\kappa(x, x') = \phi(x)^\top \phi(x')$ $\kappa(x, x')$ is in the original space! (enable computation of classifiers in high-dimensional space without performing computations in this high-dimensional space.)

**Prediction** $y = \phi(x)^\top w_*$ is Expensive, so use Rep theorem: $\phi(x)^\top w_* = \phi(x)^\top \phi(X) \alpha_* = \sum_{n=1}^{N} \kappa(x, x_n) \alpha_n$ (Non linear in X space) ($\neq$ feature space) (Linear)

**Kernels:** linear $\kappa(x, x') = x^\top x'$, $\phi(x) = x$ | Quad $\kappa(x, x') = (xx')^2$, $\phi(x) = x^2$ | Poly: $\kappa(x, x') = (x_1 x'_1 + x_2 x'_2 + x_3 x'_3)^2$, $\phi(x) = [x_1^2, x_2^2, x_3^2, \sqrt{2} x_1 x_2, \sqrt{2} x_2 x_3, \sqrt{2} x_1 x_3]^\top$ | RBF: $\kappa(x, x') = e^{-\gamma \|x - x'\|^2}$, $\phi(x) = e^{-x^2}(\dots, \frac{2^{k/2} x^k}{\sqrt{k!}} \dots)$ • $\kappa(x, x') = \ell(\kappa(x, x'))$

**Create kernels:** ⊕ lin Combo: $\kappa(x, x') = \alpha \kappa_1(x, x') + \beta \kappa_2(x, x')$ for $\alpha, \beta \geq 0$ Products: $\kappa(x, x') = \kappa_1(x, x') \kappa_2(x, x')$ • $\kappa(x, x') = e^{3\kappa_1(\kappa(x, x'))}$

**Mercers condition:** given K, ensures $\phi$ iff: Kernel fct is symmetric $\kappa(x, x') = \kappa(x', x)$ Kernel matrix psd $K = (\kappa(x_i, x_j))_{i,j=1}^N \geq 0 \forall N \geq 0 \forall (x_i)_{i=1}^N$

**SVM:** Hard-SVM: Max-margin separating HP: $\max_{w, \|w\|=1} \min_n [w^\top x_n]$ such that $\forall n: y_n x_n^\top w \geq 0$ Equiv to: $\max_{M \in \mathbb{R}, w, \|w\|=1} M$ such that $\forall n, y_n x_n^\top w \geq M$ $M = \frac{1}{\|w\|}$ (LS) (Margin) (Correctly classified)
$\Rightarrow \min_w \frac{1}{2} \|w\|^2$ such that $\forall n, y_n x_n^\top w \geq 1$ **Soft-SVM:** Not linearly separable! → use slack variable: $y_n x_n^\top w \geq 1 - \xi_n$
$\min_{w, \xi} \frac{\lambda}{2} \|w\|^2 + \frac{1}{N} \sum_{n=1}^{N} \xi_n$ s.t. $\forall n, y_n x_n^\top w \geq 1 - \xi_n$ and $\xi_n \geq 0$ equiv to: $\min_w \frac{\lambda}{2} \|w\|^2 + \frac{1}{N} \sum_{n=1}^{N} [1 - y_n x_n^\top w]_+$ → ERM for hinge loss with ridge $\lambda$ (If $y_n x_n^\top w \geq 1$, then $\xi_n = 0$ If $y_n x_n^\top w < 1$, $\xi_n = 1 - y_n x_n^\top w$)

**Margin-base losses** ($z = y x^\top w$) MSE($\eta$) = $(1 - \eta)^2$ Hinge($\eta$) = $[1 - \eta]_+$
Logistic($\eta$) = $\frac{\log(1 + \exp(-\eta))}{\log 2}$ all upper Bound for 0-1 loss

**Opt:** Convex duality: How to get $w$ for SOFT-SVM ⇒ Find $G(w, \alpha)$ so $\min_w L(w) = \min_w \max_\alpha G(w, \alpha) \geq \max_\alpha \min_w G(w, \alpha)$ (Primal) (Dual)
1) $[1 - y_n x_n^\top w]_+ = \max_{\alpha_n \in [0,1]} \alpha_n(1 - y_n x_n^\top w)$ $\min_w L(w) = \min_w \max_{\alpha \in [0,1]^N} \frac{1}{N} \sum_{n=1}^{N} \alpha_n(1 - y_n x_n^\top w) + \frac{\lambda}{2} \|w\|^2$ → G Convex in w Concave in $\alpha$ so we can invert min max
2) $\partial_w$ in $\nabla_w G(w, \alpha) = -\frac{1}{N} \sum_{n=1}^{N} \alpha_n y_n x_n + \lambda w = 0 \Rightarrow w(\alpha) = \frac{1}{\lambda N} \sum_{n=1}^{N} \alpha_n y_n x_n = \frac{1}{\lambda N} X^\top Y \alpha$ (diag(y)) 3) Max: $\min_w L(w) = \max_{\alpha \in [0,1]^N} \frac{1 - \alpha}{N} - \frac{1}{2\lambda N^2} \alpha^\top \underbrace{YXX^\top Y}_{\text{PSD}} \alpha$ (Kernel of data) ($x_n > 0$ are SV)

$\alpha_n = 0$: $x_n \checkmark$ Not margin, $\alpha_n \in (0,1)$: $x_n \checkmark$ margin, $\alpha_n = 1$: $x_n$ X

**Ethics:** We want knowledge, not stereotypes. ML can't distinguish. data might discriminate minorities. Three fairness criterias: Can't satisfy them all. Independence: $R$ independent of $A \Longrightarrow$ equal acceptance rate
Separation: $R$ independent of $A$, conditional on $Y \Longrightarrow$ equal error rate
Sufficiency: $Y$ independent of $A$ conditional on $R \Longrightarrow$ calibration by group

## Left Column

**K-NN:** $nbh_{S_{train},k}: \mathcal{X} \to \mathcal{X}^k$  $x \mapsto$ {the $k$ elements of $S_{train}$ closest to $x$}  Slow: $\Theta(N)$

For reg: $f_{S_{train},k} = \frac{1}{k}\sum_{k\in K} y_n$ (avg of points around)  For classif: $f_{S_{train},k} = $ majority $\{y_n : x_n \in nbh\}$

Small $K$: Low Bias, High Var. Use odd $K$. When $K=N$: cst prediction.

**Curse of D:** When dimension $d$ gets large (with fixed $N$), hard to find a neighboring point.  $X, Y = [0,1]^d \times \{0,1\}$

**Bayes Classifier:** min $L$ on all classif: $f_*(x) = 1_{\eta(x) \geq 1/2}$
↳ Proba min of misclassif  $\eta(x) = \mathbb{P}(Y=1|X=x)$

**Bayes Risk:** $L(f_*) = \mathbb{P}(f_*(X) \neq Y) = \mathbb{E}_{X \sim D_X}[\min\{\eta(X), 1-\eta(X)\}]$

**Bound for 1-NN:** $\mathbb{E}_{S_{train}}[L(f_{S_{train}})] \leq 2L(f_*) + 4c\sqrt{d}N^{-\frac{1}{d+1}}$  ↳ Lipschitz

If $N \to \infty$, bound is $2L(f_*)$. But we need $N \propto d^{(d+1)/2}$ to have constant error. Curse of D!

---

**CNN.** Conv: filter $f$, size $K$, stride $S$. $x_{n,m}^{(1)} = \sum_{k,l=0}^{K-1} f_{k,l} \cdot x_{nS+k,mS+l}^{(0)}$
⇒ Same filter used → weight sharing. Value depends on close values.
⇒ 0 and valid padding. Can also use multiple filters.  $W_{out} = \frac{W_{in} + 2P - F}{S} + 1$

**Conv Layer** has multiple filters. HP: size, padding, stride
**Pooling** Max or Avg → reduce spatial dim. HP: size, type, stride
ReLU after each Conv layer to make it Non-linear.
**Backpropa:** 1) Backpropa indep 2) Sum grad of edges with sam weights
**ResNet:** Add $X$ to standard Network $F(x)$ ⇒ Skip-connection $R(x) + X$
**Data augmentation:** transform $\gamma: \mathbb{R}^d \to \mathbb{R}^d$ (same label) by rotation
**Weight decay:** $\min \mathcal{L} + \frac{\lambda}{2}\sum_l \|W^{(l)}\|_F^2$ No need for Bias.
**Dropout:** Subnetwork by keeping with proba $p^{(l)}$ each node. But use whole network when testing.

---

**Adversarial ML:** Standard vs Adv risk: $R(f) = \mathbb{E}_{\mathcal{D}}[1_{f(X)\neq Y}]$  $R_\epsilon(f) = \mathbb{E}_{\mathcal{D}}[\max_{\tilde{x}, \|\tilde{x}-x\|\leq\epsilon} 1_{f(\tilde{x})\neq y}]$

**Generating adv example:** given input $(x,y)$ and $f: x \to \{-1,+1\}$ find $\tilde{x}$ st.
a) $\|\tilde{x}-x\| \leq \epsilon$  b) the model $f$ makes a mistake ⇒ $\max_{\tilde{x}, \|\tilde{x}-x\|\leq\epsilon} 1_{f(\tilde{x})\neq y}$
**Issues:** 1) $1_c$ Not continuous 2) NN pred output $\{-1,+1\}$. Solutions: 1) Use
smooth $\mathcal{L}$ 2) Consider output before classif ($f(x) = $ Sign $(g(x))$  $P(Y=1|X)$
↳ Know $g(x)$ ⇒ White-box: Solve $\max_{\tilde{x}, \|\tilde{x}-x\|\leq\epsilon} \ell(yg(\tilde{x}))$  $\nabla_x \ell'(yg(x))\nabla_x g(x)(\propto -y\nabla_x g(x))$
$\leq 0$ since classification loss are decreasing
$\ell_2: \hat{x} = x - \epsilon y \cdot \frac{\nabla_x g(x)}{\|\nabla_x g(x)\|_2}$  $\ell_\infty: \hat{x} = x - \epsilon y \cdot$ sign$(\nabla_x g(x)) \Rightarrow$ PGD multiple steps
↳ don't $K$ g(x)
**Black-box:** (a) score-based: query conti. model scores $g(x)$  but needs LOTS of queries
$\nabla_x g(x) \approx \sum_d \frac{g(x+\alpha e_i)-g(x)}{\alpha}e_i$  (b) Decision-based: query only predicted class $f(x) \in \{-1,+1\}$
2) **Transfer attacks:** Train $\tilde{f} \approx f$ on similar data ⇒ Model stealing
**Train robust models:** Train and min on adv examples and risk.
$\min_\theta \frac{1}{N}\sum_N \max_{\tilde{x}_n, \|x_n-\tilde{x}_n\|\leq\epsilon} \ell(y_n g_\theta(\tilde{x}_n))$ ⇒ approx $\hat{x}_n^* \approx \arg\max_{\tilde{x}_n} \ell(y_n g_\theta(\tilde{x}_n))$ with PGD
2) GD w.r.t $\Theta$ using $\frac{1}{N}\sum_{n=1}^N \nabla_\theta \ell(y_n g_\theta(\hat{x}_n^*))$
**Note:** robustness comes at the cost of accuracy

---

**Matrix Facto:** $X$ $(D\times N) \approx W Z^T, W \in (D\times K), Z \in (N\times K)$ user(row),$K \times N$, item(row)
⇒ $\min_{W,Z} \mathcal{L}(W,Z) = \frac{1}{2}\sum_{(d,n)\in\Omega}[x_{dn} - (WZ^T)_{dn}]^2$ ⇒ Opt. not jointly convex and not unique
$K$: # latent feat. large $K \to$ overfit  Reg: $\mathcal{L}(W,Z) + \frac{\lambda_W}{2}\|W\|_F^2 + \frac{\lambda_Z}{2}\|Z\|_F^2$
$(\|A\|_F = \sqrt{\sum_i^n\sum_j^d a_{i,j}^2})$ SGD: Derive $\nabla$ for $(d',K)$ $(W)$, $(n',K)$ $(Z)$ for fixed $(d,n)$.
$\frac{\partial}{\partial w_{d',k}} f_{d,n}(W,Z) = \begin{cases} -[x_{dn}-(WZ^T)_{dn}]z_{nk} & \text{if } d'=d \\ 0 & \text{otherwise} \end{cases}$  $\frac{\partial}{\partial z_{n',k}} f_{d,n}(W,Z) = \begin{cases} -[x_{dn}-(WZ^T)_{dn}]w_{dk} & \text{if } n'=n \\ 0 & \text{otherwise} \end{cases}$
Cost $\Theta(K)$! indep of $N, D$!
**ALS:** Assume $X$ is full, we use coord descent:  $\min \frac{1}{2}\sum_{d=1}^D\sum_{n=1}^N [x_{dn}-(WZ^T)_{dn}]^2 = \frac{1}{2}\|X-WZ^T\|_F^2$ → fix $W$ min $Z$, vice-versa
⇒ $Z^T = (W^TW + \lambda_W I_K)^{-1}W^TX$  $W^T = (Z^TZ + \lambda_W I_K)^{-1}Z^TX^T$ → optimal  Cost $\Theta(DKN+K^3)$?
⇒ But $X$ not full $\min \sum_{(n,n)\in\Omega}[x_{dn}-(WZ^T)_{dn}]^2 \to$ 1) $\nabla_W \mathcal{L}(W,Z) = 0$ 2) $\nabla_Z \mathcal{L} = 0...$

---

**Text represent:** find $w_i \to W_i \in \mathbb{R}^K$  Words→Features  Co-occurrence matrix
$n_{i,j}$ = # contexts where $w_i$ occurs with $w_j$. Very sparse $D\times N$. Need context (size) and Voc $V$
→ Matrix Facto: log of counts: $x_{dn} = \log(n_{dn}) \to X \approx WZ^T$ where $W$ (word), $Z$ (context word)
$\min_{W,Z} \mathcal{L}(W,Z) = \frac{1}{2}\sum_{(d,n)} f_{dn}[x_{dn}-(WZ^T)_{dn}]^2$  Glove: $f_{dn} = \min\{1, (n_{dn}/n_{max})^\alpha\}, \alpha \in [0,1]$
Training w SGD or ALS: $\Theta(K)$ per step  Skip-gram (word2vec): use log-reg to
separate real word pairs $(w_d, w_n)$ from fake ones. (real: appearing in context window)
**Fast text:** Supervised: $S_n = (w_1,...,w_m)$ and $x_n \in \mathbb{R}^{|V|}$ a BoW of $S_n$ and $y_n \in \{\pm 1\}$ label:
$\min_{W,Z} \mathcal{L}(W,Z) = \sum_m \ell(y_n W Z^T x_n), W \in \mathbb{R}^{1\times K}, Z \in \mathbb{R}^{K\times|V|}$, $f$ a classifier
**Supervised:** use CNN or mat fact  **Unsupervised:** adding or avg word vectors
and train, then direct unsupervised train using sentences instead of words.

---

**LLMs:** **LM:** distribution over text. $p(\text{I saw a cat on a mat}) = p(x_1,...,x_T)$  **Next-token pred:**
$P(x_T|x_1,...,x_{T-1})$. $P(x_1|x_1,...,x_{T-1})...P(x_T)$  $p(\text{mat} | \text{I saw a cat on a}) = 0.002$
**Tokenizer:** indep from model  **AR inference:** 1) Tokenize 2) Forward 3) Proba for next token
4) Sample 5) Detoken 6) repeat  2) **Data Pretraining:** A lot of low quality to learn.
**Post training:** Small of high quality to make model usable. **Distributed learning:**
LLM need multiple GPUs to run in parallel 3) **Post training:** ZS vs FS (in-context): if
we give examples to GPT, it will perform better. Also: if we give examples of
reasoning (CoT). **Supervised finetuning:** Have labelers show desired output

## Right Column

**NN: Linear:** $f_{Lin}(x) = \phi(x)^T w$ (fixed)  NN: $f_{NN}(x) = f(x)^T w$ (learned) | $L$ hidden layers $d$ inputs $K$ neurons $W^{(l)}$ matrices $b^{(l)}$ vectors

Outputs of hidden layer $l$ given by vector: $x^{(l)} = f^{(l)}(x^{(l-1)}) := \phi((W^{(l)})^T x^{(l-1)} + b^{(l)})$  $x_j^{(l)} = \phi(\sum_{i=1}^K x_i^{(l-1)}w_{i,j}^{(l)} + b_j^{(l)})$



Params: $w_{i,j}^{(l)}, b_j^{(l)}$  $\Theta(dK + K^2L)$  Barron's theorem: $\int_{|x|\leq r}(f(x)-f_n(x))^2 dx \leq \frac{(2Cr)^2}{n}$  Also use comb of ReLU as PWL
↳ Non-linearity is important

**Representation power:** We can use sigmoids to approximate functions. Cross: $(x_1,x_2) \mapsto \phi(w(x_1-a_1)) - \phi(w(x_1-b_1)) + \phi(w(x_2-a_2)) - \phi(w(x_2-b_2))$

**MSE:** $\min_{w_{i,j}^{(l)}, b_j^{(l)}} \mathcal{L}(f) = \frac{1}{2N}\sum_{n=1}^N (y_n - f(x_n))^2$ where $f$ is the fct of NN → Non-Convex! | Note: We could add a Reg term.

**SGD:** $(w_{i,j}^{(l)})_{t+1} = (w_{i,j}^{(l)})_t - \gamma\frac{\partial \mathcal{L}_n}{\partial w_{i,j}^{(l)}}$  $(b_i^{(l)})_{t+1} = (b_i^{(l)})_t - \gamma\frac{\partial \mathcal{L}_n}{\partial b_i^{(l)}}$  But doing all chain rules too expensive! Solution: Backpropagation.

**NN params:** $W^{(1)}: d\times K$, $W^{(l)}: K\times K$ for $2\leq l \leq L$, $W^{(L+1)}: K\times 1$, $b^{(l)}: K\times 1$ for $1\leq l\leq L$, $b^{(L+1)}: 1\times 1$ | $x^{(l)} = f^{(l)}(x^{(l-1)}) := \phi((W^{(l)})^T x^{(l-1)}+b^{(l)}) = \phi(z^{(l)})$
⇒ Final layer $y = f^{(L+1)}(x^{(L)}) := (W^{(L+1)})^T x^{(L)} + b^{(L+1)}$ and overall function is a compo: $y = f(x^{(0)}) = f^{(L+1)} \circ f^{(L)} \circ ... \circ f^{(2)} \circ f^{(1)}$

**Goal:** Compute for all $(i,j,l)$: $\frac{\partial \mathcal{L}_n}{\partial w_{i,j}^{(l)}}$ and $\frac{\partial \mathcal{L}_n}{\partial b_i^{(l)}}$. **FW pass:** Set $x^{(0)} = x_n$, then compute $z^{(l)}$ and $x^{(l)}$  $\Theta(K^2L)$

**BW pass:** Define $\delta_S^{(l)} = \frac{\partial \mathcal{L}_n}{\partial z_i^{(l)}} = (W^{(l+1)}\delta^{(l+1)}) \odot \phi'(z^{(l)})$ — Initialize by $\delta^{(L+1)} = z^{(L+1)} - y_n$ (MSE')  $\Theta(K^2L)$

$\frac{\partial \mathcal{L}_n}{\partial b_i^{(l)}} = \delta_S^{(l)}$  $\frac{\partial \mathcal{L}_n}{\partial w_{i,j}^{(l)}} = \delta_S^{(l)} x_i^{(l-1)}$  **Issues with GD:** Grad vanishing or exploding. Solutions are:

1) **Param initialization:** Control layerwise Var of the neurons (H0)! $z_i^{(l)} \sim \mathcal{N}(0, I_K)$  $w_i^{(l+1)} \sim \mathcal{N}(0, \sigma^2 I_K)$  $z_i^{(l+1)} = ReLU(z^{(l)})^T w_i^{(l+1)}$  $\sigma^2 = 2K$ for Var $(z_i^{(l+1)}) = 1$
2) **Batchnorm:** $\tilde{z}_n^{(l)} = (z_n^{(l)} - \mu_n^{(l)})/\sqrt{(\sigma_n^{(l)})^2 + \epsilon}$  Then $\hat{z}_n^{(l)} = \gamma^{(l)} \tilde{z}_n^{(l)} + \beta^{(l)}$  Layer norm: $\tilde{z}_n^{(l)} = (z_n^{(l)} - \mu_n^{(l)}1_K)/\sqrt{(\sigma_n^{(l)})^2 + \epsilon}$  then Same

---

**Transformers:** $f: $ seq $\to$ seq  input (T) →token→ $T$ tokens of dim $\Delta$ —$L$ transf→ $T\times b$ → output  Transformer: LN, SA (between), MLP (within), SC

**Input:** Tokenized → Embedding: map each token ID into vector $w_i \in \mathbb{R}$ → $X = [w_0, w_{123}, ... w_{07}]^T \in \mathbb{R}^{T\times b}$
$W_{emb} = [w_1 ... w_{N_{voc}}] \in \mathbb{R}^{N_{voc}\times D}$ learnt via backpropa. $X = [e_{:,1},...,e_{:,T}]^T W_{emb}$. Position: $X = [e_{:,1},...,e_{:,T}]^T W_{emb} + [e_1,...,e_T]^T W_{pos}$
↳ SA doesn't account for position

**SA:** A: tok → tok using learned index. Weighted avg. ⇒  $Z = $ softmax$\left(\frac{XW_Q W_K^T X^T}{\sqrt{D_K}}\right) XW_V$  $\Theta(T^2)$  $W_K, W_Q$ $(b\times b_K)$, $W_V$ $(b\times b_V)$
→ Can also do it for $H$ SA: $Z_h$

**MLP:** mix info within each token: $MLP(X) = \varphi(XW_1)W_2$ → applied to each token. Output: single or multiple (simple)

**Transformers** are used for Encoders (classification), Decoders (chat GPT), Encoder-Decoder (translation)

---

**Unsupervised learning:** Model-Based clustering: **K-Means:** Fix $K$ clusters and find Cluster centroids $\mu_k$ and Assignment $z_{nk} \in \{0,1\}$

**Objective function:**  $\min_{z,\mu} \mathcal{L}(z,\mu) = \sum_{n=1}^N \sum_{k=1}^K z_{nk}\|x_n - \mu_k\|_2^2$  s.t. $\mu_k \in \mathbb{R}^D, z_{nk} \in \{0,1\}, \sum_{k=1}^K z_{nk} = 1$  (WCSS) W: within

**Algo:** 1) input data $D$, # clusters $K$  2) Initialize random $\mu_K$  3) Assign each point to new center  $z_{nk} = \begin{cases} 1 & \text{if } k = \arg\min_k \|x_n-\mu_k\|_2^2 \\ 0 & \text{otherwise} \end{cases}$
4) Update centroids (means): $\mu_k = \frac{\sum_n^N z_{nk}x_n}{\sum_n^N z_{nk}}$ ⇒ repeat 3) 4) until convergence  **Iteration cost** $\Theta(NKd)$  Every iteration lowers the WCSS

**Convergence:** coordinate descent: $z^{(t+1)} = \arg\min_z \mathcal{L}(z, \mu^{(t)})$  $\mu^{(t+1)} = \arg\min_\mu \mathcal{L}(z^{(t+1)}, \mu)$ → optimum not guaranteed
$\min_{z,\mu} \mathcal{L}(z,\mu) = \|X^T - MZ^T\|_{Frob}^2$  **Challenges:** 1) Depends on initialization ⇒ K-means++: spread initial centroids $\Theta(\log K)$ optimal solution
2) Which $K$ optimal? ⇒ Elbow method of $K$ on WCSS. 3) Can't model other shapes than spheres. 4) sensitive to outliers
5) Hard assignments in overlapping clusters.

**GMM:** Calculate distrib of clusters. $K$ Gaussians. $Z$ (unobserved) multinomial → $P(z=k) = \pi_k \to \sum_k \pi_k = 1$
**MLE:** $\max_\theta \sum_{n=1}^N \log\left(\sum_{k=1}^K \pi_k \mathcal{N}(x_n|\mu_k, \Sigma_k)\right)$ but Non-Convex, no unique opt, unbounded. → Need **EM algorithm**
**Hard EM:** $\Theta^{(0)} = \{\pi_k, \mu_k, \Sigma_k\}_{k\in[K]}$. **E:** Predict most likely class for each data point: $z_n^{(t)} = \arg\max_z \pi_z \mathcal{N}(x_n|\mu_z, \Sigma_z)$  We have labeled data
**M:** MLE of $\Theta$ using $D_c = \{(x_n, z_n^{(t)})\}_{n\in[N]}$: $\theta^{(t)} = \arg\max_\theta p(D_c^{(t)}|\theta) = \arg\max_\theta \sum_{n=1}^N \log(\pi_{z_n} \mathcal{N}(x_n|\mu_{z_n}, \Sigma_{z_n}))$. K-means: Hard EM, $\pi_k = \frac{1}{K}, \Sigma_k = \sigma^2$
**Soft EM:** use posteriors $\Theta^{(0)} = \{\pi_k, \mu_k, \Sigma_k\}_{k\in[K]}$  **E:** Compute posteriors: $q_k^{(t)}(x_n) = p(z_n = k|x_n, \theta^{(t-1)}) = \frac{\pi_k \mathcal{N}(x_n|\mu_k^{(t-1)}, \Sigma_k^{(t-1)})}{\sum_{k'=1}^K \pi_{k'}\mathcal{N}(x_n|\mu_{k'}^{(t-1)}, \Sigma_{k'}^{(t-1)})}$
**M:** MLE of $\Theta$ using $q_k^{(t)}(x_n)$  $\pi_k^{(t)} = \frac{1}{N}\sum_n q_k^{(t)}(x_n)$, $\mu_k^{(t)} = \frac{\sum_{n=1}^N q_k^{(t)}(x_n)x_n}{\sum_{n=1}^N q_k^{(t)}(x_n)}$  $\Sigma_k^{(t)} = \frac{\sum_{n=1}^N q_k^{(t)}(x_n)(x_n - \mu_k^{(t)})(x_n-\mu_k^{(t)})^T}{\sum_{n=1}^N q_k^{(t)}(x_n)}$  Iteration cost $\Theta(NKD^2)$
**General EM:** Not only GMM: $\mathcal{L}(\theta) = \mathcal{L}(q,\theta) + KL(q\|p)$ where $q$ is distribution $q(z)$ s.t. $\sum_z q(z) = 1$  coordinate ascent algorithm on $\mathcal{L}(q,\theta)$
**Goal:** Max incomplete $\mathcal{L}(\theta) = \log(P(x|\theta))$ by max lower bound: $\mathcal{L}(q,\theta) = \log\left(\frac{P(x, z|\theta)}{q(z)}\right)$
**E:** Compute $P(z|x, \theta^{(t)})$: $q^{(t+1)} = \arg\max_q \mathcal{L}(q, \theta^{(t)})$  Write $\log p(X, Z|\theta)$ and calculate: $\mathbb{E}_Z[\log p(X, Z|\theta)|X, \theta^{(t)}]$
**M:** Max $\theta^{(t+1)} = \arg\max_\theta \mathbb{E}_Z[\log p(X, Z|\theta)|X, \theta^{(t)}]$ ⇒ $\theta^{(t+1)} = \arg\max_\theta \mathcal{L}(q^{(t+1)}, \theta)$  Note: $\mathcal{L}(\theta^{(t+1)}) \geq \mathcal{L}(q^{(t+1)}, \theta^{(t+1)}) \geq$

alternately computing a lower bound on the log likelihood for the current parameter values (E) and then maximizing this bound to obtain the new parameter values (M)  $\mathcal{L}(q^{(t+1)}, \theta^{(t)}) \geq \mathcal{L}(\theta^{(t)})$

**Initialization:** Weights: uniform  Means: Rand or K-Means++  Var: empirical

---

**Self-supervised:** uses input data to create label or tasks ⇒ $g: X \to (X_{in}, X_{out})$ then learn $f: X_{in} \to X_{out}$ (e.g. MLM)
**MLM:** Learn to predict hidden words. **BERT:** Finds Masked words, learn order of the sentence and classifies
**Joint embed. methods:** rather than $f: x_{in} \to x_{out}$, learn $f: X \to \mathbb{R}^d$ s.t. $f(x_n) \approx f(x_n)$  Problem: constant $f$.  Solution:
**BYOL:** $f_1(x_1) \approx f_2(x_2)$ and **Contrastive learning:** similarity of ⊕ example (pair better than w ⊖: $s(f(x), f(x^+)) > s(f(x), f(x^-))$  $f_1$ encoder, $f_2$ projector
**SimCLR:** 1. Classifies $L(x) = -\mathbb{E}\left[\log\frac{\exp(s(f(x), f(x^+)))}{\sum_{i=1}^N \exp(s(f(x), f(x_i)))}\right]$  2. Map encoder output to similarity space: $f(x) = f_2 \circ f_1(x)$
3. Cosine similarity: $s(e_1, e_2) = \frac{\langle e_1, e_2\rangle}{\|e_1\|_2\|e_2\|_2}/\gamma$  4. Data augmentation  **CLIP:** Max similarity of captions and images

---

**Gen Models:** instead of model $P(y|x)$, model $P(x)$. **Explicit:** model prob distrib of $x$  **Implicit:** generate samples according to it.
**GAN:** Generator $G(\theta)$, Discriminator $D(\varphi)$. $D$ judges what $G$ produces. $\theta^* \in \arg\min_\theta \mathcal{L}^G(\theta, \varphi^*)$  $\varphi^* \in \arg\min_\varphi \mathcal{L}^D(\theta^*, \varphi)$
1) $D$ distinguishes real vs. fake sample. 2) $G$ wants to fool $D$. 3. Obj: $\min_G \max_D \mathbb{E}_{x\sim}[\log D(x)] + \mathbb{E}_{z\sim}[\log(1-D(G(z)))]$  Losses:
$\mathcal{L}_D(G, D) = \max_D \mathbb{E}[\log D(x)] + \mathbb{E}[\log(1-D(G(z)))]$  $\mathcal{L}_G(G, D) = \min_G \mathbb{E}[\log(1-D(G(z)))]$  4. opt solution: $P_g = P_d = -\log 4$ | **CGAN:** Condition
**Diff models:** add noise and train model to recover the data. Consider a MC with $x_0 \sim P_0$ and forward transition such
that $x_{t+1} \sim P(\cdot|x_t)$. Forward decompo: $p(x_{0:T}) = p_0(x_0)\prod_{t=1}^T p(x_t|x_{t-1})$ where at each step, the marginal satisfies
$p(x_t) = \int p(x_t|x_{t-1})p(x_{t-1})dx_{t-1}$. Using bayes $(p(x_t|x_{t+1}) = p(x_{t+1}|x_t)p(x_t)/p(x_{t+1}))$ ⇒ $p(x_{0:T}) = p_T(x_T)\prod_{t=0}^{T-1} p(x_t|x_{t+1})$  Backward!

We use **ancestral sampling:** Sample $X_T \sim p_T(\cdot)$ then: for $k = T-1,...,0$: Sample $X_k \sim p(\cdot|X_{k+1})$

$P_0 = p_{data} \in \mathbb{P}(\mathbb{R}^d)$ | Forward transition: $P(x_{t+1}|x_t) = \mathcal{N}(x_{t+1}; \alpha x_t, (1-\alpha^2)\mathbb{I}_d)$ | $P(x_t|x_0) = \mathcal{N}(x_t; \alpha^t x_0, (1-\alpha^{2t})\mathbb{I}_d)$  MC converges
to $P_{ref} = \mathcal{N}(x; 0, I_d)$ for $T\to\infty$ ⇒ for $T$ large enough: $P_T(x) \approx P_{ref}(x)$. **Problem:** BW transition $P(x_t|x_{t+1})$ needs to be approx
⇒ $p(x_t|x_{t+1}) = p(x_{t+1}|x_t)\exp[\log p(x_t) - \log p(x_{t+1})] \approx \mathcal{N}(x_t; (2-\alpha)x_{t+1} + (1-\alpha^2)\nabla\log p(x_{t+1}), (1-\alpha^2)\mathbb{I}_d)$
Approximate score with a NN: $\theta^* = \arg\min_\theta \frac{1}{2}\sum_{t=1}^T \mathbb{E}_{X_0}\mathbb{E}_{X_t|X_0}\|s_\theta(t, X_t) - \nabla\log p(X_t|X_0)\|^2$  Score