

1 Convexity

- A function is convex iff a line segment between two points on the function's graph always lies above the function.
- A function $h(u)$, $u \in \mathbb{R}^D$ is convex if $\forall u, v \in \mathbb{R}^D, 0 \leq \lambda \leq 1$: $h(\lambda u + (1-\lambda)v) \leq \lambda h(u) + (1-\lambda)h(v)$ (Stictly convex if $\lambda \neq 0, 1$)
- A strictly convex function has a unique global minimum w^* . For convex functions, every local minimum is a global minimum.
- Sums of cvx funcs \Rightarrow cvx (MSE+Lin model = cvx)

2 Optimisation

- Local min $w^* \rightarrow \exists \epsilon > 0$ s.t. $L(w^*) \leq L(w) \forall w / \|w - w^*\| \leq \epsilon$
- Global minimum w^* , $L(w^*) \leq L(w) \forall w \in \mathbb{R}^D$
Gradient Descent (GD): $w^{(t+1)} := w^{(t)} - \gamma \nabla L(w^{(t)})$, $\gamma > 0$
Gradient Descent for Linear MSE: $\mathcal{O}(N \times D)$
 $L(w) := \frac{1}{N} \sum_{n=1}^N (y_n - x_n^\top w)^2 = \frac{1}{N} w^\top e$, $\nabla L(w) = -\frac{1}{N} X^\top e$
 $e := y - x$
Stochastic Gradient Descent (SGD), $\mathcal{O}(D)$
 $L(w) = \frac{1}{N} \sum_{n=1}^N \mathcal{L}_n(w)$, $w^{(t+1)} := w^{(t)} - \gamma \nabla \mathcal{L}_n(w^{(t)})$
- Cheap and unbiased estimate of the gradient $E[\nabla \mathcal{L}_n(w)] = \frac{1}{N} \sum_{n=1}^N \nabla \mathcal{L}_n(w) = \nabla(\frac{1}{N} \sum_{n=1}^N \mathcal{L}_n(w)) = \nabla \mathcal{L}(w)$
Mini-batch SGD: $w^{(t+1)} := w^{(t)} - \gamma g$
 $g := \frac{1}{|B|} \sum_{n \in B} \nabla \mathcal{L}_n(w^{(t)})$, subset $B \subseteq [N]$ of training samples

Non-Smooth Optimization

- Convexity for differentiable functions: $\mathcal{L}(u) \geq \mathcal{L}(w) + \nabla \mathcal{L}(w)^\top (u - w) \quad \forall u$
Sugradients

- A vector $g \in \mathbb{R}^D$ s.t. $\mathcal{L}(u) \geq \mathcal{L}(w) + g^\top (u - w) \quad \forall u$ is called a subgradient of the function \mathcal{L} at w .

This definition makes sense for objectives \mathcal{L} which are not necessarily differentiable (and not even necessarily convex).

- L convex and differentiable at $w \rightarrow$ only subgradient $g = \nabla \mathcal{L}(w)$.

Subgradient Descent

$w^{(t+1)} := w^t - \gamma g$ for g a subgradient to L

Convex Sets

- Intersections of convex sets are convex
- Functions onto convex sets are unique
- Euclidean definition: $P_S(w) := \{v \in \mathbb{R}^D \mid \|v - w\| \leq \epsilon\}$

Projected Gradient Descent SGD + Projection step
Idea: add a projection onto C after every step:
 $P_C(w^{(t)}) := \arg \min_w \|w - w^{(t)}\|$, $w^{(t+1)} := P_C(w^{(t)} - \gamma \nabla \mathcal{L}(w^{(t)})]$

Constrained + Unconstrained Problems

$\min_{w \in C} \mathcal{L}(w) \rightarrow \min_w \mathcal{L}(w) + P(w)$

Implementation Issues

Optimality: 1st order: if L is cvx and $\nabla \mathcal{L}(w^*) = 0 \rightarrow$ global 2nd order: L potentially non-convex, if $\nabla^2 \mathcal{L}(w^*) = 0$, $\nabla^2 \mathcal{L}(w^*) \geq 0 \rightarrow$ w local minimum ($\nabla^2 \mathcal{L}(w) := \frac{\partial^2 \mathcal{L}}{\partial w^2}(w)$ expensive)

3 Least Squares (Linear regression + MSE)

$L(w) = \frac{1}{N} \sum_{n=1}^N (y_n - x_n^\top w)^2 = \frac{1}{N} (y - Xw)^\top (y - Xw)$,
- Proof of convexity:

1) Composition of a linear function with a convex function (the square function).
2) Verify definition, for any $\lambda \in [0, 1]$ and w, w' : $\mathcal{L}(\lambda w + (1 - \lambda)w') - (\lambda \mathcal{L}(w) + (1 - \lambda)\mathcal{L}(w')) \geq 0$.

LHS = $-\frac{1}{2N} \lambda (1 - \lambda) \|X(w - w')\|_2^2 < 0$

3) Hessian positive semidefinite (all its eigenvalues are non-negative):
 $H(w) = \frac{1}{N} \nabla^2 \mathcal{L}(y - Xw)^\top (y - Xw) = \frac{1}{N} X^\top X$

Singular value decomposition (PSD): $X = USV^\top$ U and V are orthogonal matrices, and S is a diagonal matrix with the singular values s_i on the diagonal.

- Now find its minimum $\nabla \mathcal{L}(w) = -\frac{1}{N} X^\top (y - Xw) = 0$

Closed form (if Gram matrix $X^\top X \in \mathbb{R}^{D \times D}$ is invertible)

$w^* = (X^\top X)^{-1} X^\top y$

Invertibility and Uniqueness

- $X^\top X \in \mathbb{R}^{D \times D}$ invertible iff $\text{rank}(X) = D$.

- Proof: assume $\text{rank}(X) < D \rightarrow \exists u \neq 0$ s.t. $Xu = 0 \rightarrow X^\top Xu = 0 \rightarrow \text{rank}(X^\top) < D \rightarrow X^\top$ is not invertible.

Rank Deficiency and Ill-Conditioning

- In practice, X is often rank deficient.
- If $D > N$, we always have $\text{rank}(X) < D$ (row rank = col. rank)

- If $D \leq N$, but some of the columns x are (nearly) collinear, then the matrix is illconditioned \rightarrow numerical issues

4 Maximum Likelihood

Gaussian distribution and Independence

$p(y | \mu, \sigma^2) = \mathcal{N}(y | \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp[-\frac{(y-\mu)^2}{2\sigma^2}]$, Σ psd.

$\mathcal{N}(\mu | \mu, \Sigma) = \frac{1}{\sqrt{(2\pi)^D \det(\Sigma)}} \exp[-\frac{1}{2} (y - \mu)^\top \Sigma^{-1} (y - \mu)]$

A probabilistic model for least-squares

- $y_n = x_n^\top w + \epsilon_n$ where the ϵ_n zero mean Gaussian RV
 $p(y | X, w) = \prod_{n=1}^N p(y_n | x_n, w) = \prod_{n=1}^N \mathcal{N}(y_n | x_n^\top w, \sigma^2)$

Defining cost with log-likelihood

$L_{LL}(w) := \log p(y | X, w) = -\frac{1}{N} \sum_{n=1}^N (y_n - x_n^\top w)^2 + \text{const.}$

Maximum-likelihood estimator (MLE)

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

Properties of MLE

- MLE is a sample approximation to the expected log-likelihood: $\mathcal{L}_{LL}(w) \approx \mathbb{E}_{p(y | X, w)} [\log p(y | X, w)]$

- MLE is consistent, i.e., it will give us the correct model assuming that we have a sufficient amount of data. $w_{MLE} \xrightarrow{p} w_{true}$

- The MLE is asymptotically normal, i.e.,

$(w_{MLE} - w_{true}) \xrightarrow{d} \frac{1}{N} \mathcal{N}(\mathcal{L}_{MLE}, \mathbf{F}^{-1}(w_{true}))$

where $\mathbf{F}(w) = -\mathbb{E}_{p(y | X, w)} [\frac{\partial \log p(y | X, w)}{\partial w}]$ is the Fisher information.

- MLE is efficient, i.e., it achieves the Cramer-Rao lower bound.

Covariance of $w_{MLE} = \mathbf{F}^{-1}(w)$.

Laplace distribution: $p(y_n | x_n, w) = \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2} |y_n - x_n^\top w|}$

5 Regularization: $\min_w \mathcal{L}(w) + \Omega(w)$

Ridge Regression: Regularization + \mathcal{L}_{MSE}

- Euclidean norm (L_2 norm): $\Omega(w) = \|w\|_2^2$, $\|w\|_2^2 = \sum_i w_i^2$, $\min_w \frac{1}{2N} \sum_{n=1}^N (y_n - x_n^\top w)^2 + \lambda \|w\|_2^2 = \min_w \frac{1}{N} \|y - Xw\|^2 + \lambda \|w\|_2^2$

Explicit solution for w:

$\nabla \mathcal{L} = \nabla \mathcal{L}_{MSE} + \lambda \nabla \Omega$

$\nabla \mathcal{L}_{MSE} = (X^\top X + \lambda I)^{-1} X^\top y$, $\lambda = \frac{2}{N} \lambda$

Ridge Regression Fights Ill-Conditioning

- Lifting the eigenvalues : The eigenvalues of $(X^\top X + \lambda I)$ are all at least λ and so the inverse always exists.

- Proof: Use the Eigenvalue decomposition of $X^\top X$ as USU^\top , $X^\top X + \lambda I = USU^\top + \lambda I = U(S + \lambda I)U^\top$

L1-Regularization: The Lasso (L1-norm + \mathcal{L}_{MSE})

$\min_w \frac{1}{N} \sum_{n=1}^N |y_n - x_n^\top w|^2 + \lambda \|w\|_1$ where $\|w\|_1 := \sum_i |w_i|$

6 Model Selection

Generalization Error

- Expected error over all samples drawn from distribution \mathcal{D} : $L_\phi(f) = \mathbb{E}_{(x,y) \sim \mathcal{D}} [\ell(y, f(x))]$

Empirical Error

- approximate the true error by averaging the loss function over the dataset: $L_S(f) = \frac{1}{|S|} \sum_{(x_n, y_n) \in S} \ell(y_n, f(x_n))$, (RV)

- unbiased estimator of the true error

- Law of large number: $L_S(f) \xrightarrow{|S| \rightarrow \infty} L_\phi(f)$ but fluctuations!

Generalization gap: $|L_\phi(f) - L_S(f)|$

- Claim: given a model f and a test set S_{test} $\sim \mathcal{D}$ i.i.d. (not used to learn f) and a loss $\ell(\cdot, \cdot) \in [a, b]$:

MLE for logistic regression

- Assumption: The inputs \mathbf{x} do not depend on the parameter w we choose:

$\mathcal{L}(w) = p(\mathbf{y} | \mathbf{X}, w) = p(\mathbf{X} | w)p(\mathbf{y} | \mathbf{X}, w)$

$\mathcal{L}(w) = \prod_{n=1}^N p(x_n | w) p(y_n | x_n, w) = \prod_{n=1}^N p(y_n | x_n, w)$

$\Pi_{n=1}^N p(y_n | x_n, w) = \prod_{n=1}^N (1 - p(x_n | w))^{1-y_n} p(x_n | w)^{y_n}$

$\mathcal{L}(w) \propto \prod_{n=1}^N \sigma(x_n^\top w)^{y_n} [1 - \sigma(x_n^\top w)]^{1-y_n}$

Minimum of the Negative Log Likelihood (NLL)

$-\log(p(\mathbf{y} | \mathbf{X}, w)) = -\log \left(\prod_{n=1}^N \sigma(x_n^\top w)^{y_n} [1 - \sigma(x_n^\top w)]^{1-y_n} \right)$

$= -\sum_{n=1}^N -y_n x_n^\top w + \log(1 + e^{x_n^\top w})$

- True loss: $E[\mathcal{L}(y_n, f(x_n))] = L_\phi(f)$

Hoeffding's inequality: a simple concentration bound

- Claim: Let $\Theta_1, \dots, \Theta_N$ be a sequence of i.i.d. random variables with mean $\mathbb{E}[\Theta]$ and range $[a, b]$. Then, for any $\epsilon > 0$

$\Pr[\frac{1}{N} \sum_{n=1}^N \Theta_n - \mathbb{E}[\Theta] \geq \epsilon] \leq 2e^{-2N\epsilon^2/(b-a)^2}$

- Concentration bound: the empirical mean is concentrated around its mean

A. Use it with $\Theta_n = \ell(y_n, f(x_n))$ B. Equating $\delta = 2e^{-2N\epsilon^2/(b-a)^2}$

we get $\epsilon = \sqrt{\frac{(b-a)^2 \ln(2/\delta)}{2N}}$

How far is each of the K test errors $L_{test}(f_k)$ from the true $L_\phi(f_k)$?

- Claim: we can bound the maximum deviation for all K candidates, by $\Pr[\max_k |L_\phi(f_k) - L_{test}(f_k)| \geq \epsilon] \leq O(\epsilon^{-1})$

- The error decreases as $O(1/\sqrt{|S_{test}|})$ with the number of test points

- test K hyper-parameters, the error only goes up by $\sqrt{\ln(K)}$

\Rightarrow can test many different models without incurring a large penalty

- It can be extended to infinitely many models

Proof: A simple union bound

- Newton's method minimizes the quadratic approximation:

$L(w) = \nabla L(w)^\top \nabla^2 L(w) (w - w_0) + \frac{1}{2} w^\top \nabla^2 L(w) w$

$\nabla^2 L(w) = \frac{1}{N} \sum_{n=1}^N \sigma(x_n^\top w) (1 - \sigma(x_n^\top w)) x_n x_n^\top$

$\nabla^2 L(w) = \frac{1}{N} \mathbf{X}^\top \mathbf{X} w$

Convexity of the loss function L

- The Hessian $\nabla^2 L(w) = \frac{1}{N} \sum_{n=1}^N \sigma(x_n^\top w) (1 - \sigma(x_n^\top w)) x_n x_n^\top$, is psd

$\nabla^2 L(w) = \frac{1}{N} \sum_{n=1}^N \sigma(x_n^\top w) (1 - \sigma(x_n^\top w)) x_n x_n^\top \geq 0$

$\Rightarrow L$ is convex since $\nabla^2 L(w) \geq 0$

Newton's method uses second order information

- Newton's method approximates the quadratic approximation

$L(w) = \nabla L(w)^\top \nabla^2 L(w) (w - w_0) + \frac{1}{2} w^\top \nabla^2 L(w) w$

$\nabla^2 L(w) = \frac{1}{N} \sum_{n=1}^N \sigma(x_n^\top w) (1 - \sigma(x_n^\top w)) x_n x_n^\top$

$\nabla^2 L(w) = \frac{1}{N} \mathbf{X}^\top \mathbf{X} w$

Newton's method uses second order information

- Newton's method minimizes the quadratic approximation

$L(w) = \nabla L(w)^\top \nabla^2 L(w) (w - w_0) + \frac{1}{2} w^\top \nabla^2 L(w) w$

$\nabla^2 L(w) = \frac{1}{N} \sum_{n=1}^N \sigma(x_n^\top w) (1 - \sigma(x_n^\top w)) x_n x_n^\top$

$\nabla^2 L(w) = \frac{1}{N} \mathbf{X}^\top \mathbf{X} w$

Newton's method is the gradient descent method

- A decomposition in three terms

Interested in the expectation of the true risk over training set S

$E_{\mathcal{D}}[\mathcal{L}(f_S)] = E_{\mathcal{D}}[\mathbb{E}_{\mathcal{D}-S}[\ell(f(x_0) - f_S(x_0))^2]] \leq \sum_n \Pr[\ell(f(x_0) - f_S(x_0)) \geq \epsilon] \leq \Pr[\ell(f(x_0) - f_S(x_0)) \geq \epsilon]$

$\leq 2K e^{-2N\epsilon^2/(b-a)^2}$

- $\delta = 2Ke^{-2N\epsilon^2/(b-a)^2} \rightarrow \epsilon = \sqrt{\frac{(b-a)^2 \ln(2K/\delta)}{2N}}$

- Let $k^* = \arg \min_k L_{test}(f_k)$ (smallest true risk) and $\hat{k} = \arg \min_k L_{test}(f_k)$ (smallest empirical risk) then:

$\Pr[L_{test}(f_{k^*}) \geq L_{test}(f_{\hat{k}})] \leq \Pr[L_{test}(f_{k^*}) \geq L_{test}(f_{\hat{k}}) + \epsilon]$

$\Pr[L_{test}(f_{k^*}) \geq L_{test}(f_{\hat{k}}) + \epsilon] \leq \Pr[\ell(f_{k^*} - f_{\hat{k}}) \geq \epsilon]$

$\Pr[\ell(f_{k^*} - f_{\hat{k}}) \geq \epsilon] = \Pr[\ell(f_{k^*} - f_{\hat{k}}) \geq \epsilon] \leq \Pr[\ell(f_{k^*} - f_{\hat{k}}) \geq \epsilon]$

$\Pr[\ell(f_{k^*} - f_{\hat{k}}) \geq \epsilon] \leq 2e^{-2N\epsilon^2/(b-a)^2}$

$\Pr[\ell(f_{k^*} - f_{\hat{k}}) \geq \epsilon] \leq 2e^{-2N\epsilon^2/(b-a)^2}$

Soft SVM: when training set is not linearly separable

- Maximize margin but allow some constraints to be violated

- Introduce positive slack variables ξ_1, \dots, ξ_N and replace the constraints with $y_n - w^\top x_n \geq 1 - \xi_n$

$\Pr[\min_{w, \xi} \frac{1}{2} \|w\|^2 + \frac{1}{N} \sum_{n=1}^N \xi_n \leq M \text{ s.t. } \forall n, y_n - w^\top x_n \geq 1 - \xi_n \text{ and } \xi_n \geq 0]$

$\Pr[\min_{w, \xi} \frac{1}{2} \|w\|^2 + \frac{1}{N} \sum_{n=1}^N \xi_n \leq M \text{ s.t. } \forall n, y_n - w^\top x_n \geq 1 - \xi_n \text{ and } \xi_n \geq 0]$

$\Pr[\min_{w, \xi} \frac{1}{2} \|w\|^2 + \frac{1}{N} \sum_{n=1}^N \xi_n \leq M \text{ s.t. } \forall n, y_n - w^\top x_n \geq 1 - \xi_n \text{ and } \xi_n \geq 0]$

$\Pr[\min_{w, \xi} \frac{1}{2} \|w\|^2 + \frac{1}{N} \sum_{n=1}^N \xi_n \leq M \text{ s.t. } \forall n, y_n - w^\top x_n \geq 1 - \xi_n \text{ and } \xi_n \geq 0]$

$\Pr[\min_{w, \xi} \frac{1}{2} \|w\|^2 + \frac{1}{N} \sum_{n=1}^N \xi_n \leq M \text{ s.t. } \forall n, y_n - w^\top x_n \geq 1 - \xi_n \text{ and } \xi_n \geq 0]$

$\Pr[\min_{w, \xi} \frac{1}{2} \|w\|^2 + \frac{1}{N} \sum_{n=1}^N \xi_n \leq M \text{ s.t. } \forall n, y_n - w^\top x_n \geq 1 - \xi_n \text{ and } \xi_n \geq 0]$

$\Pr[\min_{w, \xi} \frac{1}{2} \|w\|^2 + \frac{1}{N} \sum_{n=1}^N \xi_n \leq M \text{ s.t. } \forall n, y_n - w^\top x_n \geq 1 - \xi_n \text{ and } \xi_n \geq 0]$

$\Pr[\min_{w, \xi} \frac{1}{2} \|w\|^2 + \frac{1}{N} \sum_{n=1}^N \xi_n \leq M \text{ s.t. } \forall n, y_n - w^\top x_n \geq 1 - \xi_n \text{ and } \xi_n \geq 0]$

$\Pr[\min_{w, \xi} \frac{1}{2} \|w\|^2 + \frac{1}{N} \sum_{n=1}^N \xi_n \leq M \text{ s.t. } \forall n, y_n - w^\top x_n \geq 1 - \xi_n \text{ and } \xi_n \geq 0]$

$\Pr[\min_{w, \xi} \frac{1}{2} \|w\|^2 + \frac{1}{N} \sum_{n=1}^N \xi_n \leq M \text{ s.t. } \forall n, y_n - w^\top x_n \geq 1 - \xi_n \text{ and } \xi_n \geq 0]$

$\Pr[\min_{w, \xi} \frac{1}{2} \|w\|^2 + \frac{1}{N} \sum_{n=1}^N \xi_n \leq M \text{ s.t. } \forall n, y_n - w^\top x_n \geq 1 - \xi_n \text{ and } \xi_n \geq 0]$

$\Pr[\$

Batch Normalization

$\bar{z}_n^{(t)} = \frac{z_n^{(t)} - \mu_n^{(t)}}{\sqrt{(z_n^{(t)})^2 + \epsilon}}$ (Component-wise) where $\mu_n^{(t)} = \frac{1}{M} \sum_{k=1}^M z_n^{(t)}$ and $\epsilon = \sqrt{(\sigma_n^{(t)})^2 + \epsilon}$

$(\sigma_n^{(t)})^2 = \frac{1}{M} \sum_{k=1}^M (z_n^{(t)} - \mu_n^{(t)})^2$, and $\epsilon \in \mathbb{R}_{\geq 0}$ is a small value added for numerical stability

- Introduce learnable parameters $\gamma^{(t)}, \beta^{(t)} \in \mathbb{R}^K$ to reverse the normalization: $\bar{z}_n^{(t)} = \gamma^{(t)} \odot \bar{z}_n^{(t)} + \beta^{(t)}$
- Scale-invariance: BN ($a \odot \bar{z}_n^{(t)} + b$) = BN ($\bar{z}_n^{(t)}$)
- Inference: Estimate $\hat{\mu}^{(t)} = \mathbf{E}[\mu_n^{(t)}]$ and $\hat{\sigma}^{(t)} = \mathbf{E}[(\sigma_n^{(t)})^2]$ during training, use these for inference
- Layer Normalization**

$\bar{z}_n^{(t)} = \frac{z_n^{(t)} - \mu_n^{(t)} - 1_K}{\sqrt{\sum_{k=1}^K (\mu_n^{(t)} - \mu_n^{(t)})^2}}$, where $\mu_n^{(t)} = \frac{1}{K} \sum_{k=1}^K z_n^{(t)}(k)$ and $(\sigma_n^{(t)})^2 = \frac{1}{K} \sum_{k=1}^K (z_n^{(t)}(k) - \mu_n^{(t)})^2$, and $\epsilon \in \mathbb{R}_{\geq 0}$

- Learnable parameters $\gamma^{(t)}, \beta^{(t)} \in \mathbb{R}^K$: $\bar{z}_n^{(t)} = \gamma^{(t)} \odot \bar{z}_n^{(t)} + \beta^{(t)}$
- Normalize across features, independently for each observation
- common alternative, widely used for transformers and text data
- No batch dependency, use the same for training and inference

15 Convolutional Networks

$x_n^{(t)} = \sum_k f_{k,t} \cdot x_{n-k+1}^{(0)}$ (f is the learnable filter)

- Same filter at every position - weight sharing
- Translation equivariance: shifted input results in shifted output

Training:

1. Run backpropagation as if the weights were not shared
2. Sum the gradients of all edges that share the same weight

Weight Decay: l_2 -regularization

- Regularize weights without bias: $\min \mathcal{L} + \frac{\lambda}{2} \sum_i \|\mathbf{W}^{(i)}\|_F^2$

- Favors small w which can aid in generalization and opti

$(w_{i,j}^{(t)})_{t+1} = (w_{i,j}^{(t)})_t - \eta \nabla \mathcal{L} - \eta \lambda (w_{i,j}^{(t)})_t = (1 - \eta \lambda) (w_{i,j}^{(t)})_t - \eta \nabla \mathcal{L}$

weight decay

- Interaction with BatchNorm:

- BN($\mathbf{W}\mathbf{X}$) = BN($\alpha\mathbf{W}\mathbf{X}$) for $\alpha \in \mathbb{R}_{>0}$ (assuming $\epsilon \approx 0$)

16 Transformers

- Self-Attention (SA): mixes information between tokens

- Multi-Layer Percep. (MLP): mixes information within each token

- Skip connections are widely used

- Layer normalization (LN) placed at the start of a residual branch

Text Token Embeddings

- Tokenization: split text into sequence of tokens (predefined)

- Convert each token ID $i \in \{1, \dots, N_{\text{vocab}}\}$ into $\mathbf{w}_i \in \mathbb{R}^D$

- Matrix multiplication $\mathbf{W} \cdot \mathbf{e}_i = \mathbf{W}_{:,i} = \mathbf{w}_i$ (with $\mathbf{W} \in \mathbb{R}^{D \times N_{\text{vocab}}}$)

- \mathbf{W} learned via backpropagation

- Input sequence of T tokens leads to an input matrix $\mathbf{X} \in \mathbb{R}^{T \times D}$

Attention learning input-dependent weighted average

- Input tokens: $V \in \mathbb{R}_{\geq 0}^{T \times D}$, Output tokens: $Z = P \cdot V$, Weighting coefficients $P \in [0, 1]^{T \times T}$ valid probability distributions over input $\sum_{j=1}^T p_{ij} = 1$

- Query tokens: $Q \in \mathbb{R}_{\geq 0}^{T \times D_K}$, Key tokens: $K \in \mathbb{R}_{\geq 0}^{T \times D_K}$

- Determine weight $p_{i,j}$ based on how similar \mathbf{q}_i and \mathbf{k}_j are

- Use inner product to obtain raw similarity scores.

- Normalize with softmax (scaled temperature by $\sqrt{D_K}$) to obtain a probability distribution (applied on each row independently)

- Scaling \rightarrow uniformity at initialization and faster convergence

Self-Attention

- V, K, Q all derived from same input token sequence $X \in \mathbb{R}^{T \times D}$


- Values : $V = XW_V \in \mathbb{R}^{T \times D}, W_V \in \mathbb{R}^{D \times D}$

- Keys : $K = XW_K \in \mathbb{R}^{T \times D_K}, W_K \in \mathbb{R}^{D \times D_K}$

- Queries : $Q = XW_Q \in \mathbb{R}^{T \times D_K}, W_Q \in \mathbb{R}^{D \times D_K}$

- W_Q, W_V, W_K are learned params.

- $Z = \text{softmax} \left(\frac{XW_Q^T X^T}{\sqrt{D_K}} \right) XW_V$



Multi-Head Self-Attention

- Run H parallel "heads" in parallel $Z_h = \text{softmax} \left(\frac{XW_Q h W_V^T X^T}{\sqrt{D_K}} \right) XW_h \in \mathbb{R}^{T \times D_h}, W_{V,h} \in \mathbb{R}^{D \times D_h}, W_{K,h} \in \mathbb{R}^{D \times D_h}, W_{Q,h} \in \mathbb{R}^{D \times D_h}, Z = [Z_1, \dots, Z_H]W_o$ where $W_o \in \mathbb{R}^{HD_h \times D}$ is learned via backprop

Positional Information

- Attention by itself does not account for the order of input

- Positional encoding in the network $pos : \{1, \dots, T\} \rightarrow \mathbb{R}^D$

- e.g., W_{pos} corresponding to each token's position t to the input embedding. $W_{pos} \in \mathbb{R}^{D \times T}$ is learned via backprop

MLP: Mixing Information within Tokens

- Apply the same transformation to each token independently: $MLP(X) = \phi(XW_h)W_o, W_1, W_2 \in \mathbb{R}^{D \times D}$ learned via backprop

- Single output: transfo, to a special task-specific input token or to the average tks. Multi output: transfo to each token indep.

Vision Transformer Architecture

- The receptive field is the whole image after one SA layer

- ViTs require more data than CNNs, reduced inductive bias in extracting local features

- Model attends to image regions semantically relevant for clf

Encoders & Decoders

- Encoders: fixed output size and process all inputs simultaneously

- Decoders: Auto-regressively sample the next token as $x_{t+1} \sim \text{softmax}(f(x_1, \dots, x_t))$

17 Adversarial ML

- Standard risk: average zero-one loss over X : $R(f) = \mathbb{E}_{\{Y=1\} \times \mathcal{X}} [\mathbb{P}_f(Y=1) \neq Y]$ i.e. minimise proba of wrong pred.

- Adversarial risk: average zero-one loss over small, worst-case perturbations of X : $R_e(f) = \mathbb{E}_\theta [\max_{\|\tilde{x}-x\| \leq \epsilon} \mathbb{P}_f(\tilde{x}) \neq Y]$

Generating adversarial examples

- Task: given an input (x, y) and a model $f: \mathcal{X} \rightarrow \{-1, 1\}$ find an input \tilde{x} such that $\|\tilde{x} - x\| \leq \epsilon$ the model f makes a mistake

- Trivial case: x already misclassified \rightarrow no action required
- General case: find \tilde{x} such that $\text{atf}(\tilde{x}) \neq y$ and $\|\tilde{x} - x\| \leq \epsilon$ i.e. $\tilde{x} \in B_r(x) \cap \{f(\tilde{x}) \neq y\}$
- Optimization problem with respect to the inputs

- Problem: optimizing the indicator function is difficult: 1) The indicator function 1 is not continuous 2) The NN prediction f outputs discrete class values $\{-1, 1\}$

- Replace the difficult problem involving the indicator with a smooth problem $\max_{\|\tilde{x}-x\| \leq \epsilon} 1_{\{f(\tilde{x}) \neq Y\}} \rightarrow \max_{\|\tilde{x}-x\| \leq \epsilon} \ell(y(\tilde{x}))$

- decreasing, margin-based (i.e., dependent on $y * g(x)$) clf loss

White-Box attacks

- Solve $\max_{\|\tilde{x}-x\| \leq \epsilon} \ell(y(\tilde{x}))$ knowing g

- $\nabla_x \ell(y(\tilde{x})) \odot \nabla_x g(x)$: $\nabla_x \ell(y(\tilde{x}))$ with $y(\ell'(y(\tilde{x}))) \leq 0$ since classification losses are decreasing

- Move in direction of $\propto -\nabla_x g(x)$

- Interpretation $f(x) = \text{sign}(y(g(x)))$: If $y = 1$ we want to decrease $g(x)$ and follow $-\nabla_x g(x)$. If $y = -1$ we want to increase $g(x)$ and follow $\nabla_x g(x)$

- Use ℓ not $y(g(\tilde{x}))$: extend to multi-class clf and robust training.

- linearize the loss $\ell(\tilde{x}) := \ell(y(g(\tilde{x})))$

$\max_{\|\tilde{x}-x\| \leq \epsilon} \ell(\tilde{x}) \approx \max_{\|\tilde{x}-x\| \leq \epsilon} \ell(x) + \nabla_x \ell(x) \cdot (\tilde{x} - x) = \ell(x) + \max_{\|\tilde{x}-x\| \leq \epsilon} \ell'(\tilde{x})(\tilde{x} - x) = \ell(x) + \max_{\|\tilde{x}-x\| \leq \epsilon} \ell'(\tilde{x}) \cdot \nabla_x g(x)$

- Max inner product under a norm constraint,

- Simple problem for which we can get a closed-form solution depending on the norm used to measure the perturb size $\|\delta\|$

Two-step attack

- These updates can be done iteratively and combined with a projection step on the feasible set (i.e., balls $\ell_2 \leq \ell_\infty$ here)

- Projected Gradient Descent (PGD) attack

- δ norm: $\delta^{(t+1)} = \Pi_{B_r(x)}[\delta^t + \alpha \cdot \frac{\nabla \ell(x+\delta^t)}{\|\nabla \ell(x+\delta^t)\|}]$, $\Pi_{B_r(x)}(\delta) = \begin{cases} \delta & \text{if } \|\delta\|_2 \leq \epsilon \\ \epsilon \cdot \frac{\delta}{\|\delta\|_2} & \text{otherwise} \end{cases}$

- ℓ_∞ norm: $\delta^{(t+1)} = \Pi_{B_r(x)}[\delta^t + \alpha \cdot \text{sign}(\nabla \ell(x+\delta^t))]$, $\Pi_{B_r(x)}(\delta)_i = \begin{cases} \epsilon & \text{if } |\delta_i| \geq \epsilon \\ \epsilon \cdot \text{sign}(\delta_i) & \text{otherwise} \end{cases}$

- Multi-step attack

- These updates can be done iteratively and combined with a projection step on the feasible set (i.e., balls $\ell_2 \leq \ell_\infty$ here)

- Projected Gradient Descent (PGD) attack

- δ norm: $\delta^{(t+1)} = \Pi_{B_r(x)}[\delta^t + \alpha \cdot \frac{\nabla \ell(x+\delta^t)}{\|\nabla \ell(x+\delta^t)\|}]$, $\Pi_{B_r(x)}(\delta) = \begin{cases} \delta & \text{if } \|\delta\|_2 \leq \epsilon \\ \epsilon \cdot \frac{\delta}{\|\delta\|_2} & \text{otherwise} \end{cases}$

- ℓ_∞ norm: $\delta^{(t+1)} = \Pi_{B_r(x)}[\delta^t + \alpha \cdot \text{sign}(\nabla \ell(x+\delta^t))]$, $\Pi_{B_r(x)}(\delta)_i = \begin{cases} \epsilon & \text{if } |\delta_i| \geq \epsilon \\ \epsilon \cdot \text{sign}(\delta_i) & \text{otherwise} \end{cases}$

- Marginalized z_n out

- The gradients with backprop w.r.t. inputs, not parameters!

- Black-box attacks $g(x)$ unknown

- Obtaining a surrogate model, costly + no guarantee of success

- Query-based methods often require a lot of queries (10k-100k), Query-based gradient estimation

- Score-based: we can query the continuous model scores $g(x) \in \mathbb{R}$.

$\nabla_x g(x) \approx \sum_{i=1}^d \frac{g(x+a_i) - g(x)}{a_i}$

- Decision-based: we can query only the predicted class $f(x) \in \{-1, 1\}$, similar techniques can be adapted for the decision-based case

Transfer Attacks

- Train a similar surrogate model $\hat{f} \approx f$ on similar data

- Model stealing (query f given some unlabeled inputs $(x_i, f(x_i))_{i=1}^N$) can facilitate transfer attacks.

Adversarial training

- Adversarial training: the goal is to minimize the adversarial risk:

$R_e(f) = \mathbb{E}_{\{Y=1\} \times \mathcal{X}} [\max_{\|\tilde{x}-x\| \leq \epsilon} \ell_f(\tilde{x})]$

- Unknown \rightarrow approximate it with a sample average + classification loss is non-continuous \rightarrow use a smooth loss \Rightarrow

$\min_{\theta} \sum_{i=1}^N \max_{\|\tilde{x}_i-x_i\| \leq \epsilon} \ell_f(y_i, g(\tilde{x}_i))$

1) $\forall x_i, \tilde{x}_i \approx \arg \max_{\|\tilde{x}_i-x_i\| \leq \epsilon} \ell_f(y_i, g(\tilde{x}_i))$

2) GD step w.r.t. θ using $\frac{1}{N} \sum_{i=1}^N \nabla_\theta \ell_f(y_i, g(\tilde{x}_i))$

- Advantages:

- Maximize the lower bound w.r.t. θ

- More interpretable gradients

- fully compatible with SGD

- Disadvantages

- Increased comp time: prop to the number of PGD steps

- Robustness-accuracy tradeoff: too large ϵ = worse accuracy

18 Fairness criterions in classification

Use an algorithm to produce a score function $R = r(X)$ (given)

Sensitive attributes

- No fairness through unawareness: removing/ignoring sensitive attribute is not solving the problem

- Many features slightly correlated with the sensitive attribute can be used to recover the attribute

Three fundamental fairness criteria

Independence: $A \perp R \mid Y$, Separation: $A \perp R \mid Y$, Sufficiency: $A \perp Y \mid R$ any of these three criteria are mutually exclusive

Independence: equal acceptance rate

$P(D=1 \mid A=a) = P(D=1 \mid A=b) \quad (\text{equal FP})$

$P(D=0 \mid A=Y-1, a=A) = P(D=0 \mid Y=1, A=b) \quad (\text{equal FN})$

- Self-reliance: $P(Y=1 \mid R=r, A=a) = P(Y=1 \mid R=r, A=b)$

- Calibration: $P(Y=1 \mid R=r, A=a) = r$, $P(Y=0 \mid R=r, A=a) = 1-r$, sufficiency How to achieve fairness criteria

- Post-processing: adjust your learned classifier so that it becomes uncorrelated with the sensitive attribute A

- At training time: add regularization

- Pre-processing: adjust your features so that they become uncorrelated with the sensitive attribute A

19 Clustering

- Clusters are groups of points whose inter-point distances are small compared to the distances outside the cluster.

- Find "prototypic" points $\mu_1, \mu_2, \dots, \mu_K$ and cluster assignments $z_n \in \{1, 2, \dots, K\}$ for all $n = 1, 2, \dots, N$, data vectors $\mathbf{x}_n \in \mathbb{R}^D$.

K-means clustering

$\text{Assume } K \text{ is known, } \min_{\theta, \mu} \mathcal{L}(\mathbf{x}, \mu) = \sum_{n=1}^N \sum_{k=1}^K z_{nk} \|\mathbf{x}_n - \mu_k\|_2^2$

s.t. $\mu_k = \frac{1}{\sum_{n=1}^N z_{nk}} \sum_{n=1}^N z_{nk} \mathbf{x}_n$

- This cost is not jointly convex w.r.t. W and Z , nor identifiable as $(w^*, z^*) \Leftrightarrow (Bw^*, -B^{-1}z^*)$

22 Matrix Factorization

Given items (movies) $d = 1, 2, \dots, D$ and users $n = 1, 2, \dots, N$, $X \in \mathbb{R}^{D \times N}$ (sparse)

Algorithm

$\mathbf{X} \approx WZ^T, W \in \mathbb{R}^{D \times K}, Z \in \mathbb{R}^{N \times K}$ tall matrices $K \ll N, D$

$\min_{\theta, \mu} \mathcal{L}(\mathbf{x}, \mu) = \frac{1}{2} \sum_{(d,n) \in \mathcal{E}} \|\mathbf{x}_d - (WZ^T)_d\|^2$

- This cost is not jointly convex w.r.t. W and Z , nor identifiable as $(w^*, z^*) \Leftrightarrow (Bw^*, -B^{-1}z^*)$

Generative adversarial networks GANs

- generator G , params θ and discriminator D params φ .

$\theta^* = \arg \min_{\theta} \mathcal{L}^G(\theta, \varphi^*)$ $\varphi^* = \arg \min_{\varphi} \mathcal{L}^D(\theta^*, \varphi)$

- $\mathcal{L}^G = -\mathcal{L}^D$ zero-sum game and is a minimax problem.

1. The discriminator D distinguishes real vs. fake samples: p_z "noise" distr., e.g. $\mathcal{N}(0, I)$, p_g real data distr., $D: x \mapsto y \in \mathcal{Y}$

2. The generator aims at fooling the discriminator that its samples are real: $G: z \mapsto x$, such that if $z \approx p_z$, then hopefully $x \approx p_g$ the "fake" data distribution

3. Objective: $\min_G \max_{\varphi} \mathbb{E}_{z \sim p_z} [\log D(x)] + \mathbb{E}_{x \sim p_g} [\log(1 - D(G(x)))]$

$\mathcal{L}_D(G, D) = \mathbb{E}_z [\log(1 - D(G(z)))]$

$\mathcal{L}_G(G) = \min_G \mathbb{E}_x [\log(1 - D(G(x)))]$

4. Theoretical Solution: the optimum is reached when $p_g = p_d$ and the optimal value is $\log 4$

Alternating Least Squares

$\sum_{k=1}^D \sum_{n=1}^N \max_{\|\tilde{x}_n\| \leq 1} [\mathbf{x}_n - (WZ^T)_n]^2 = \frac{1}{2} \|\mathbf{X} - \mathbf{WZ}^T\|_{\text{Frob}}^2$

$Z^{\dagger} = (W^T W + \lambda_{\text{reg}})^{-1} W^T \mathbf{X}, W^{\dagger} := (Z^{\dagger} Z + \lambda_{\text{reg}})^{-1} Z^{\dagger} \mathbf{X}$

- Cost: needed to invert a $K \times M$ matrix

23 Text Representation

- For