# Linear Regression

Linear Regression: MSE (mean-sq error) objective function: $J_D(\vec{w},b) = \frac{1}{N}\sum_{i=1}^{N}\left(y^{(i)} - (\vec{w}^T\vec{x}^{(i)}+b)\right)^2$

- gradient descent: $\vec{g} \leftarrow \nabla_\theta J(\theta) = \frac{1}{N}\sum_{i=1}^{N}(\vec{\theta}^T\vec{x}^{(i)} - y^{(i)})\vec{x}^{(i)}$, $\vec{\theta} \leftarrow \vec{\theta} - \gamma\vec{g}$   $O(MN)$ per iteration.   [$\to J^{(i)}(\theta)$]

  ↳ local optimisation algo ⟺ converge to local min if exists ⟺ GDLR globally convergent. (soln exists, may not be unique)

minimising MSE: $\nabla J(\theta)=0 \Leftrightarrow X^TX\theta = X^Ty \Leftrightarrow$ lstsq(X,y) ⟸ find argmin $\|\theta\|$ s.t. normal eqn satisfied

↑ soln unique ⟺ N (# examples) ≥ # of LI features (e.g. 3D → 1 or 2 LI features if on line/plane)

# Stochastic Gradient Descent

Stochastic Gradient Descent: $\vec{\theta} \leftarrow \vec{\theta} - \gamma\nabla_\theta J^{(i)}(\vec{\theta})$, $i\sim U[1,N]$ or shuffle $[1,N]$

- epoch = single pass through training data (GD ⟺ 1 update/epoch, SGD ⟺ N updates)

| | #steps to convergence | compute/step |
|---|---|---|
| GD | $O(\log 1/\epsilon)^*$ | $O(NM)$ |
| SGD | $O(1/\epsilon)^*$ | $O(M)$ |

* WHP under assumptions  
↓learning rate ⟹ SGD behaves like GD  
both initial training MSE large due to uninformed init.

MLE: $\theta^{MLE} = \arg\max_\theta \prod_{i=1}^{N} p(y^{(i)}|\vec{x}^{(i)},\vec{\theta}) = \arg\max_\theta p(D|\theta)$, $\theta^{MAP} = \arg\max_\theta p(D|\theta)p(\theta)$

[recall] (online) perceptron: $\hat{y} = \text{sign}(\vec{\theta}^T\vec{x}^{(t)})$, if misclassified, $\vec{\theta} \leftarrow \vec{\theta} + y^{(t)}\vec{x}^{(t)}$

# Logistic Regression

Logistic Regression: $p(y|\vec{x},\vec{\theta}) = y=1 \Leftrightarrow \sigma(\vec{\theta}^T\vec{x})$, $y=0 \Leftrightarrow 1-\sigma(\vec{\theta}^T\vec{x})$ ← $\sum_{i=1}^{N}(y^i\ln\sigma^i + (1-y^i)\ln(1-\sigma^i))$

- $\ell(\vec{\theta}) = \sum_{i=1}^{N}\log p(y^{(i)}|\vec{x}^{(i)},\vec{\theta})$, $J(\vec{\theta}) = -\frac{\ell(\vec{\theta})}{N} \Rightarrow \frac{\partial J^{(i)}}{\partial\vec{\theta}} = -(y^{(i)} - \sigma(\vec{\theta}^T\vec{x}))\vec{x}^{(i)}$

- prediction: $y=1 \Leftrightarrow \sigma(\vec{\theta}^T\vec{x}) \geq 0.5$  ← $\vec{\theta} += \gamma(y^{(i)} - \sigma(\vec{\theta}^T\vec{x}))\vec{x}^{(i)}$ always updated.

# Regularization

Regularization: prevent overfitting (idea: Occam's razor): $\hat{\theta} = \arg\min_\theta J(\theta) + \lambda r(\theta)$

L1/Lasso: $\|\vec{\theta}\|_1 = \sum_{m=1}^{M}|\theta_m|$, L2/Ridge: $\|\vec{\theta}\|_2 = \sum_{m=1}^{M}\theta_m^2$

→ fit model params. → true function classifier → training error used to choose $h\in\mathcal{H}$  * validation error to choose $\mathcal{H}$ (hyperparameters)

# Learning Theory

Learning Theory: true error rate $R(h) = E_{x\sim p^*}[\mathbb{I}(c^*(x)\neq h(x))]$ is unknown

empirical risk/training error $\hat{R}(h) = E_{x\sim D}[?] = \frac{1}{M}\sum_{i=1}^{M}\mathbb{I}[y^{(i)}\neq h(x^{(i)})]$

↑ $c^*$ may be unachievable, best achievable (true) risk minimizer $h^* = \arg\min_{h\in\mathcal{H}} R(h)$ unknown  
only know empirical risk minimizer $\hat{h} = \arg\min_{h\in\mathcal{H}} \hat{R}(h)$. overfitting $= \hat{R}(h) - R(h)$

PAC (Probably Approximately Correct) criterion: $P(|R(h) - \hat{R}(h)| \leq \epsilon) \geq 1-\delta$ $\forall h\in\mathcal{H}$

↑ $\epsilon$ is diff between true & empirical risk, $\delta$ is probability of "failure"

→ realizable if $c^*\in\mathcal{H}$, agnostic if $c^*$ might/might not be in $\mathcal{H}$

ERM (empirical risk minimisation) on $\mathcal{H}$ with $M$ training ex.

→ sample complexity = M needed in order to satisfy PAC for given $\epsilon,\delta$ → finite if $|\mathcal{H}|<\infty$, infinite if $|\mathcal{H}|=\infty$.

→ Bayesian view: $\theta$ is RV, described by prior & posterior, MAP  [find estimator that → $\theta$ quickly, even for worst-case distrib of data]  
Frequentist: $\theta$ as a constant, (regularized) MLE, consistency/convergence rates/robustness

↳ regularized MLE = MAP if regularizer = log(prior), but justification different.  
→ MAP estimate = mode of posterior; β-distribution $\propto \theta^{\alpha-1}(1-\theta)^{\beta-1}$ → shape if ↑α,β  
  → α more mass to θ=1, β→θ=0

---

$\vec{x} \leftrightarrow 6\times1$, $\vec{x}^{(i)}$ w bias (1 added)  ← weights α is $3\times(6+1)$

$\vec{a} = \vec{\alpha}\vec{x}^{(i)} \leftrightarrow a_j = \alpha_{j,b} + \sum_{i=1}^{6}\alpha_{j,i}x_i$ $\forall 1\leq j\leq 3$  ← # of neurons in hidden layer is 3

$\vec{z} = \text{ReLU}(\vec{a}) = \max(0,\vec{a})$, $\vec{z}$ is $3\times1$  ← alternative sigmoid(u) $= \frac{1}{1+e^{-u}}$

$\vec{\beta}$ $4\times(3+1)$ → prepend 1 to 1st entry of $\vec{z}$  $\frac{\partial z}{\partial a} = \mathbb{I}[a>0]$  $\sigma'(u) = \sigma(u)(1-\sigma(u))$

$\vec{b} = \vec{\beta}\vec{z} \leftrightarrow 4\times1$  * if β 1-row of non-zero values, $\vec{b}$ reducible to layer $\vec{w}$  [single neuron (ie. all d-1 neurons have output 0)]

$\hat{y} = \text{softmax}(\vec{b}) = \frac{\exp(b_k)}{\sum_{l=1}^{4}\exp(b_l)}$  $\frac{\partial\hat{y}_l}{\partial b_k} = \hat{y}_l(\mathbb{I}[k=l]-\hat{y}_k)$  $\frac{\partial l}{\partial b_k} = \sum_l \frac{\partial l}{\partial\hat{y}_l}\frac{\partial\hat{y}_l}{\partial b_k} = \hat{y}_k - y_k$  [1-hot vector]

$\ell(\vec{y},\vec{\hat{y}}) = -\sum_{i=1}^{4}y_i\ln(\hat{y}_i)$ → $\frac{\partial l}{\partial\hat{y}_i} = -\frac{y_i}{\hat{y}_i}$ & $\sum_l y_l = 1$

→ LINEAR: $\vec{x} - \boxed{\vec{w}} - \vec{z}$  ($\frac{\partial l}{\partial\beta_{kj}} = \frac{\partial l}{\partial b_k}z_j$)  used to calculate $\frac{\partial l}{\partial\alpha}$  (no bias bc. $\beta_{k,0}$ are not affected by values of α)  
fwd: return $\vec{w} @ \vec{x}^{(i)}$ ← cache $\vec{x}^{(i)}$  
bck: $\frac{\partial l}{\partial\vec{w}} = \frac{\partial l}{\partial\vec{z}}\vec{x}^T$ (np.outer(dz, $\vec{x}^{(i)}$)), return $\frac{\partial l}{\partial\vec{x}} = (\vec{w}^*)^T\frac{\partial l}{\partial\vec{z}}$ → step: $\vec{w} -= lr * \frac{\partial l}{\partial\vec{w}}$

→ ReLU: fwd cache $\max(0,\vec{x})$, bck $\partial[in] = \partial[out] * (\text{cache} > 0)$  
  Sigmoid: fwd cache $1/(1+\exp(-x))$, bck $\partial[in] = \partial[out]$ cache $(1-\text{cache})$

→ Softmax Cross Entropy: fwd return $(\hat{y}, -\ln\hat{y}[y])$, bck $y_k = [0,\cdots,1,\cdots,0]$ for $\frac{\partial l}{\partial b_k}$  [cross entropy loss, k-index]  
  ↳ if separate: $g_b = g_{\hat{y}}(\text{diag}(\hat{y}) - \hat{y}\hat{y}^T) + l = -\vec{y}^T\ln\hat{y}$, $g_{\hat{y}} = -\frac{\vec{y}}{\hat{y}} \cdot g_l$ ← 1

① for finite $\mathcal{H}$ s.t. $c^*\in\mathcal{H}$ (realizable), and arbitrary distribution $p^*$, if # labelled training data pts  
$M \geq \frac{1}{\epsilon}\left(\ln(|\mathcal{H}|) + \ln(\frac{1}{\delta})\right)$ then with prob $\geq 1-\delta$, $h\in\mathcal{H}$ with $\hat{R}(h)=0$ have $R(h)\leq\epsilon$

Pf: $E$ = event $\exists h\in\mathcal{H}$ with $\hat{R}(h)=0$, $R(h)>\epsilon$, then ⟹ $P(E)<\delta$  ← if $\exists h$ with $R(h)>\epsilon$  * union bound $+ \ln(1+?)\leq ?$  
$P(E) < k(1-\epsilon)^M \leq |\mathcal{H}|(1-\epsilon)^M \Rightarrow \ln P(E) < \ln|\mathcal{H}| + M\ln(1-\epsilon)$  Ex: $\mathcal{H}$ = conjunctions over d boolean variables ⟹ $|\mathcal{H}| = 3^d$ [1,0, absent]

↳ Cor: given training data set S s.t. $|S|=M$, all $h\in\mathcal{H}$ with $\hat{R}(h)=0$ have $R(h) \leq \frac{1}{M}(\ln|\mathcal{H}| + \ln(\frac{1}{\delta}))$

② for finite $\mathcal{H}$ and arbitrary dist. $p^*$, if # labeled training dpts satisfies $R(h)<\epsilon$ ↗ w.p. at least $1-\delta$  
$M \geq \frac{1}{2\epsilon^2}\left(\ln|\mathcal{H}| + \ln(\frac{2}{\delta})\right)$, then wp at least $1-\delta$, all $h\in\mathcal{H}$ satisfy $|R(h)-\hat{R}(h)|\leq\epsilon$  $R(h) + \frac{1}{M}(\ln|\mathcal{H}|)$

↳ Cor: ... s.t. $|S|=M$, all $h\in\mathcal{H}$ have $R(h) \leq \hat{R}(h) + \sqrt{\frac{1}{2M}(\ln|\mathcal{H}| + \ln\frac{2}{\delta})}$ w.p. at least $1-\delta$.  [actual train error + regularizer]

Def: $\mathcal{H}$ shatters set of pts if it can classify them all possible ways  
Sauer's lemma: sps $S_H(M) = 2^d$ for $M\leq d$, but $S_H(d+1)<2^{d+1}$, then $S_H(M)\in O(...)$  
↳ VC($\mathcal{H}$) = size of largest set $\mathcal{H}$ can shatter ($\exists d$ pts... $\nexists d+1$ pts...)  
  → halfspaces in d dimensions: VC = d+1 ↳ so $|\mathcal{H}| \to S_H(M)$ (thms)

if $S_H(M) = 2^M$ $\forall M$, not learnable/can memo at any |D|, if bounded, can memo at non d.

**Algorithm 1** SGD

1: Initialize $\theta^{(0)}$
2:
3:
4: $s = 0$
5: **for** $t = 1, 2, \ldots, T$ **do**
6:     **for** $i \in \text{shuffle}(1, \ldots, N)$ **do**
7:         Select the next training point $(x_i, y_i)$
8:         Compute the gradient $g^{(s)} = \nabla J_i(\theta^{(s-1)})$
9:         Update parameters $\theta^{(s)} = \theta^{(s-1)} - \eta g^{(s)}$
10:         Increment time step $s = s + 1$
11:     Evaluate average training loss $J(\theta) = \frac{1}{n} \sum_{i=1}^{n} J_i(\theta)$
12: **return** $\theta^{(s)}$

$-\ell = -\log p(w^{(i)})$

**Algorithm 1** Mini-Batch SGD

1: Initialize $\theta^{(0)}$
2: Divide examples $\{1, \ldots, N\}$ randomly into batches $\{I_1, \ldots, I_B\}$
3: where $\bigcup_{b=1}^{B} I_b = \{1, \ldots, N\}$ and $\bigcap_{b=1}^{B} I_b = \emptyset$
4: $s = 0$
5: **for** $t = 1, 2, \ldots, T$ **do**
6:     **for** $b = 1, 2, \ldots, B$ **do**
7:         Select the next batch $I_b$, where $m = |I_b|$
8:         Compute the gradient $g^{(s)} = \frac{1}{m} \sum_{i \in I_b} \nabla J_i(\theta^{(s)})$
9:         Update parameters $\theta^{(s)} = \theta^{(s-1)} - \eta g^{(s)}$
10:         Increment time step $s = s + 1$
11:     Evaluate average training loss $J(\theta) = \frac{1}{n} \sum_{i=1}^{n} J_i(\theta)$
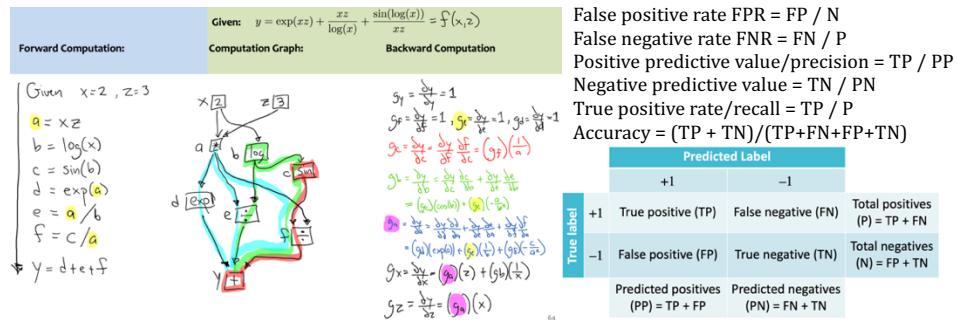12: **return** $\theta^{(s)}$

**Convexity**: $f \colon \mathbb{R}^D \to \mathbb{R}$ is convex if $\forall \mathbf{x}^{(1)}, \mathbf{x}^{(2)} \in \mathbb{R}^D$ and $0 \le c < 1$, $f\left(c\mathbf{x}^{(1)} + (1-c)\mathbf{x}^{(2)}\right) \le cf\left(\mathbf{x}^{(1)}\right) + (1-c)f\left(\mathbf{x}^{(2)}\right)$, linear functions are convex, $ax^2$ concave for $a < 0$

- MSE, MAE are convex, strictly convex iff $\mathbf{x}$ has full column rank (otherwise infinite minimizers in nullspace; adding strictly convex regulariser makes both strictly convex
- Convex → converge to global minima which might not exist (e.g. $e^x$) (exist for MSE/MAE)

<span style="color:red">**Don't use test error in making model selection decisions!**</span>

Matrix multiplication for $[M \times N][N \times P] \in O(MNP)$, matrix inverse in $O(N^3)$

**Conditional likelihood**: iid samples $D = \left\{x^{(i)}, y^{(i)}\right\}$ from a pair of RVs (unlike likelihood function with single $X$ RV with pmf $p(x \mid \theta)$), $Y$ discrete with pmf $p(y \mid x, \theta)$



False positive rate FPR = FP / N
False negative rate FNR = FN / P
Positive predictive value/precision = TP / PP
Negative predictive value = TN / PN
True positive rate/recall = TP / P
Accuracy = (TP + TN)/(TP+FN+FP+TN)

| | Predicted Label | |
|---|---|---|
| | +1 | −1 | |
| **+1** (True label) | True positive (TP) | False negative (FN) | Total positives (P) = TP + FN |
| **−1** | False positive (FP) | True negative (TN) | Total negatives (N) = FP + TN |
| | Predicted positives (PP) = TP + FP | Predicted negatives (PN) = FN + TN | |

**Achieving fairness**: (1) pre-processing data, (2) additional constraints during training, (3) post-processing predictions – premise for 1+2: if def of fairness satisfied in training data, then most models will preserve that relationship. $A$ protected label, $X$ applicant data, $Y$ pred

1. **Independence** (selection rate parity): $h(X, A) \perp A$, prop of accepted applicants same for all genders (adjust penalty for predicting +ve in class till we get parity/use diff threshold), permits laziness (alw pred +1)/susceptible to adversaries (admit some randomly)
   - Prediction rate is the same across values of $A$, $P(h = 1 \mid A = a_1) = P(h = 1 \mid A = a_2)$
2. **Separation** (FPR = FNR): $h(X, A) \perp A \mid Y$, all good/bad applicants accepted with same prob rgdless of $A$, perpetuate existing bias (only access to target var thru historical data)
   - Among individuals of the same true label, $P(h = 1 \mid Y, A)$, classifier independent of $A$
3. **Sufficiency** (PPV = NPV): $Y \perp A \mid h(X, A)$, among people who receive the same prediction, the actual probability of being positive is the same across groups, $P(Y = 1 \mid h, A)$ independent of $A$ (i.e., the info contained in $h(X, A)$ is sufficient, $A$ becomes irrelevant)

If baseline rates of label across both values of A equal, then $Y \perp A$, both S's can be achieved.