

SUPERVISED LEARNING

- + LINEAR REGRESSION
- + LINEAR CLASSIFICATION (PERCEPTRON)
- + LOGISTIC REGRESSION
- + NEAREST NEIGHBOUR
- + DECISION TREES
- + MULTICLASS SOFTMAX
- + NAIVE BAYES
- + GENERATIVE MULTICLASS GAUSSIAN.
- + SUPPORT VECTOR MACHINES

LINEAR REGRESSION

- + INPUT DATA: $(x^1, y^1), (x^2, y^2), \dots, (x^n, y^n)$
where $x^i = (x_1, x_2, \dots, x_d)$. and y^i = dependent variable
 $x^i \in \mathbb{R}^d$ and $y^i \in \mathbb{R}$ [NOTE, CONTINUOUS OUTPUT]
- + OUTPUT MECHANISM: + In a linear regression, we have some LOSS FUNCTION
e.g. $L_2(y, \hat{y}) = (y - \hat{y})^2$, $L_1(y, \hat{y}) = |y - \hat{y}|$ and a
PREDICTION MECHANISM e.g. $f_w(x) = w_0 + w_1 x_1 + \dots + w_d x_d$
where $f_w(x) = \hat{y}$.
- + We can use BASIS FNS $\phi(x_1), \dots, \phi(x_d)$ b/c linear in
params not inputs.
- + Need to MINIMIZE THE LOSS, which gives us \hat{y} closest to y
as possible. Can do this analytically. Goal is to find
 w^* minimizing L_1/L_2 LOSS FN IN VARIABLE OF w .
This is "training" linreg.
- + When done, w 's optimized. Can output \hat{y} closest to y . Model
of data.

+ KEY MATHS/ALGOS:

- + Analytical derivation of w^* minimizing loss:

Define $\hat{y} = Xw$ ($n \times d \times d \times 1 = n \times 1$ preds).

Then, the L_2 loss is: $L_2(w) = \|y - Xw\|^2$

$$\text{Then, } L_2(w) = (y - Xw)^T (y - Xw) = (y^T - (Xw)^T)(y - Xw) = (y^T - w^T X^T)(y - Xw)$$

$$\begin{aligned} \text{Then, } (y^T - w^T X^T)(y - Xw) &= y^T y - y^T Xw - w^T X^T y + w^T X^T Xw \\ &= y^T y - w^T X^T y - w^T X^T y + w^T X^T Xw \quad \} \text{ equiv.} \\ &= y^T y - 2w^T X^T y + w^T X^T Xw. \end{aligned}$$

Then, $\partial_w L(w) = -2X^T y + 2X^T Xw \rightarrow$ set to 0: $0 = 2X^T Xw - 2X^T y$

$\rightarrow w = (X^T X)^{-1} X^T y$. NOTE. IF features NOT LIN DEP, $(X^T X)^{-1}$ exists.
(use solvers: conj. grad, Cholesky decomp)

With REGULARIZATION: $w = (X^T X + \lambda I)^{-1} X^T y$

+ KEY INFORMATION:

- + KEY INFO:
 - + Requires initially designing good feats, outside of scope of algo. Can be difficult!
 - + L₂ LOSS: $L_2(y, \hat{y}) = (y - \hat{y})^2$
Optimal prediction is the CONDITIONAL MEAN $E[y|x]$
Assumes GAUSSIAN NOISE
 - + L₁ LOSS: $L_1(y, \hat{y}) = |y - \hat{y}|$
Optimal prediction is the CONDITIONAL MEDIAN.
Assumes LAPLACE NOISE
 - + L₁ LOSS VS L₂ LOSS:

L ₂ LOSS	L ₁ LOSS
+ Not resistant to data outliers (low robustness, punishes mistakes)	+ Is robust
+ Stable soln	+ Unstable soln (for a small horizontal mmt of data pt, reg. line may jump a LOT)
+ Always one soln	+ Multiple solns possible - (euclid vs. taxicab)
- + L₂ REG: $R_2(w) = \frac{\lambda}{2} \|w\|^2 = \frac{\lambda}{2} w^T w = \frac{\lambda}{2} \sum w_j^2 \sim \text{GAUSSIAN PRIOR}$
 - + Bias NOT regularized
 - + Find good $\lambda \rightarrow$ VALIDATION / CV
- + L₁ REG: $R_1(w) = \lambda \|w\| = \lambda \sum |w_j| \sim \text{LAPLACE PRIOR}$.
- + L₁ REG VS L₂ REG

L ₂ REG	L ₁ REG	
+ Analytic soln: comp. efficient.	+ Comp. inefficient b/c non-sparsity	
+ Non-sparse outs	+ Sparse outs	- "only few matrix entries non-zero"
+ No feat. sel.	+ Feat. sel.	- b/c L ₁ \rightarrow sparse coefficients
- + PROBABILISTIC VIEW: (L₂)

Assume $p(y_i|x_i) = w^T x_i + \epsilon_i, \epsilon_i \sim N(0, \sigma^2)$

$$p(y|x) = N(w^T x, \sigma^2) = \frac{\exp\left(-\frac{\|y - w^T x\|^2}{2\sigma^2}\right)}{\sqrt{2\pi\sigma^2}}$$

We know w parameterizes a distribution. Which distro to pick? Maximize LIKELIHOOD that we get y 's given x 's!

$$\begin{aligned} p(y|x) &= p(y^1|x^1)p(y^2|x^2)\dots p(y^N|x^N). \\ \log(p(y|x)) &= \log(p(y^1|x^1)p(y^2|x^2)\dots p(y^N|x^N)) \\ &= \sum_N \log p(y^i|x^i) = \sum_N \log \left(\frac{\exp\left(-\frac{\|y - w^T x\|^2}{2\sigma^2}\right)}{\sqrt{2\pi\sigma^2}} \right) = \\ &= \sum_N \left[-\frac{\|y - w^T x\|^2}{2\sigma^2} - \frac{1}{2} \log(2\pi\sigma^2) \right] \end{aligned}$$
- + MAP estimation:

$$w_{\text{MAP}} = \arg \max P(w|y, x) = p(y|w, x)p(w).$$

CLASSIFICATION

LINEARLY SEPARABLE:

- + NON-PERFECT SEPARATION: Data of two classes COMPLETELY SEPARABLE BY SINGLE LINE.
- + Simple model
- + Noise
- + Errors in data target (mislabels)
- + Simple features
- + Different feature parameterization.
- + ZERO-ONE LOSS
 - + 1 if wrong, 0 if right.
- + SURROGATE LOSS: PIECEWISE SMOOTH + CONVEX.
 - + $\ell(y, \hat{y}) \leq \hat{\ell}(y, \hat{y}) \rightarrow \text{minimizing } \hat{\ell} \text{ minimizes } \ell.$
- + METRICS:
 - + PRECISION, RECALL, ROC (PP/TP)

LOGISTIC REGRESSION

- + INPUT: (x_i^j, y_i) , $y \in \{0, 1\}$.
- + We classify with a linear decision boundary: $\sigma(z) = \frac{1}{1+e^{-z}}$. Therefore, if $\sigma(z) \geq 0.5 \rightarrow \hat{y}=1$, $\sigma(z) < 0.5 \rightarrow \hat{y}=0$.
- + To learn the weights, we have a probabilistic model.
- Need to use MAX LIKELIHOOD.

Want to maximize the $p(y|x; w)$. (likelihood). $\max_w L(w) = \max_w \prod_{i=1}^N p(y_i|x_i; w)$

$$\begin{aligned} p(y|x; w) &= p(y^1, \dots, y^N | x^1, \dots, x^N; w) \\ &= p(y^1|x^1; w)p(y^2|x^2; w)\cdots p(y^N|x^N; w). \quad (\text{IID}) \\ &= [p(y=1|x^1; w)^{y^1} p(y=0|x^1; w)^{1-y^1}] \cdots [p(y=1|x^N; w)^{y^N} p(y=0|x^N; w)^{1-y^N}] \end{aligned}$$

1) Maximize log-likelihood (neg \rightarrow minimize)

$$L(w) = \prod_{i=1}^N p(y_i|x_i; w) = \log L(w) = \log \left(\prod_{i=1}^N p(y_i|x_i; w) \right) = \sum_{i=1}^N \log(p(y_i|x_i; w)).$$

$$\text{Then, } \log L(w) = \sum_{i=1}^N \log(p(y_i|x_i; w)) = \sum_{i=1}^N \log \left[p(y=1|x^i; w)^{y^i} p(y=0|x^i; w)^{1-y^i} \right]$$

$$- \log L(w) = - \sum_{i=1}^N \log \left[p(y=1|x^i; w)^{y^i} (1-p(y=1|x^i; w))^{1-y^i} \right]$$

$$= - \left(\sum_{i=1}^N [y^i \log p(y=1|x^i; w) + (1-y^i) \log p(y=0|x^i; w)] \right)$$

$$= - \sum_{i=1}^N y^i \log p(y=1|x^i; w) - \sum_{i=1}^N (1-y^i) \log p(y=0|x^i; w)$$

$$p(y=1|x^i; w) = (1 + \exp(-w^T x^i))^{-1} \quad p(y=0|x^i; w) = 1 - p(y=1|x^i; w)$$

$$\begin{aligned}
-\log L(w) &= -\sum_{i=1}^N y_i \log((1+\exp(-z))^{-1}) \rightarrow \sum_{i=1}^N (1-y_i) \log(\exp(-z)(1+\exp(-z))^{-1}) \\
&= \sum_{i=1}^N y_i \log(1+\exp(-z)) - \sum_{i=1}^N (1-y_i) [\log(\exp(-z)) + \log((1+\exp(-z))^{-1})] \\
&= \sum_{i=1}^N y_i \log(1+\exp(-z)) - \sum_{i=1}^N (1-y_i)(-z) - (1-y_i) \log(1+\exp(-z)) \\
&= \sum_{i=1}^N y_i \log(1+\exp(-z)) + \sum_{i=1}^N (1-y_i)z + \sum_{i=1}^N (1-y_i) \log(1+\exp(-z)) \\
&= \cancel{\sum_{i=1}^N y_i \log(1+\exp(-z))} + \sum_{i=1}^N (1-y_i)z + \sum_{i=1}^N \log(1+\exp(-z)) - \cancel{\sum_{i=1}^N y_i \log(1+\exp(-z))} \\
&= \sum_{i=1}^N (1-y_i)z + \sum_{i=1}^N \log(1+\exp(-z)).
\end{aligned}$$

Now, taking the derivative:

$$\begin{aligned}
\frac{\partial L}{\partial w_j} &= \frac{\partial L}{\partial z} \frac{\partial z}{\partial w_j} = \sum_{i=1}^N (1-y_i)x_j \cdot \sum_{i=1}^N \frac{1}{1+\exp(-z)} \cdot \exp(-z) \cdot x_j = \\
&= \sum_{i=1}^N (1-y_i)x_j - \frac{x_j \exp(-z)}{1+\exp(-z)} = \sum_{i=1}^N x_j \left(1-y_i - \frac{\exp(-z)}{1+\exp(-z)} \right) = \sum_{i=1}^N x_j \left(\frac{1}{1+\exp(-z)} - y_i \right)
\end{aligned}$$

Thus,

$$\frac{\partial L}{\partial w_j} = \sum_{i=1}^N x_j (p(y=1|x_i; w) - y_i)$$

+ LOG-LOSS is more robust to outliers than L_2 !

+ MAP estimation: regularization to avoid overfitting:

$$\text{MAP: } \max_w \log \left[p(w) \prod_{i=1}^N p(y_i|x_i; w) \right]$$

$$p(w) \sim N(\mu, \sigma^2 I) \rightarrow \text{equiv to } \frac{\lambda}{2} \|w\|^2 = \frac{\lambda}{2} \sum_{j=1}^n w_j^2. \quad L_2 \text{ regularization}$$

$$+ \text{New gradient: } w_j^{(t+1)} \leftarrow w_j^{(t)} - \lambda \left(\frac{\partial L(w)}{\partial w_j} + \alpha w_j \right) = w_j^{(t)} - \lambda \frac{\partial L(w)}{\partial w_j} - \lambda \alpha w_j$$

+ TUNING HYPERPARAMS:

+ Training and validation set.

+ CAN overfit on validation set

+ LEAVE-p-OUT: $\binom{n}{p}$.

+ K-FOLD: Divide train set into k equal subsamples

k iterations: k-1 training, kth validation.

k results can be averaged for single hyperparam estimate

- PROS:

- class pred prob. view
- good acc.
- quick train
- resistant to overfit
- convex loss (log loss)
- ($N \geq 10 \cdot d$)
- Fast classify
- Model coeff = feature importance

+ CONS:

- Linear decision boundary
- Simple model, assumptions of cond. probs

1ST NEIGHBOUR

NONPARAMETRIC: while log reg is a probabilistic classifier, this is a EUCLIDEAN classifier. GROWING # OF PARAMS. NONLINEAR DECISION BOUNDARIES

+ ALGORITHM:

1. Take point x_i
2. Find k points $x^{*1} \dots x^{*k}$ s.t. each is minimal distance to x_i .
3. Vote on class. $y^{*1} \dots y^{*k}$. Assign $y_i = \text{mode}(y^{*1} \dots y^{*k})$.

+ CHOOSING k :

Generally $k < \text{sqrt}(n)$. Choose k via cross validation.

+ kNN: ISSUES + REMEDIES

- + Attributes w/ larger ranges get weighed more
↳ NORMALIZE
- + Irrelevant, correlated attributes add noise
↳ ELIM. ATTRIBUTES
- + Nonmetric attributes
↳ HAMMING DIST
- + High Dim Data
↳ DIM. REDUCTION
- + Expensive at Test Time
 - ↳ USE DIM. SUBSETS
 - ↳ PRESORT TRAINING EXAMPLES
 - ↳ APPROX DIST ONLY
 - ↳ REMOVE REDUNDANT DATA

DECISION TREES

- + Each node is a feature, "SPLIT" the universe of values of these features to reach other features. Classification \rightarrow tree descent.
LEAVES ARE CLASS LABELS.

+ ENTROPY EQUATIONS:

$$H(X) = - \sum_{x \in X} p(x) \log_2 p(x).$$

$$H(X, Y) = - \sum_{x \in X} \sum_{y \in Y} p(x, y) \log_2 p(x, y).$$

$$H(Y|X=x) = \sum_{y \in Y} p(y|x) \log_2 p(y|x)$$

$$H(Y|X) = \sum_{x \in X} p(x) H(Y|X=x) = - \sum_{x \in X} \sum_{y \in Y} p(x, y) \log_2 p(y|x)$$

$$IG(Y|X) = H(Y) - H(Y|X)$$

+ IF X, Y INDEP:

$$H(X|Y) = H(Y).$$

+ CHAIN RULE:

$$\begin{aligned} H(X, Y) &= H(X|Y) + H(Y) \\ &= H(Y|X) + H(X). \end{aligned}$$

+ DECISION TREE ALGORITHM

- 1) calculate $H(Y)$ where $Y = \text{class targets}$
- 2) calculate conditional entropy of each feature:
 $H(Y|X_1), H(Y|X_2), \dots, H(Y|X_d)$.
- 3) calculate information gain for each feature:
 $IG(Y, X_1), IG(Y, X_2), \dots, IG(Y, X_d)$.
- 4) choose attribute w/ LARGEST IG as DECISION NODE. Assume here its $X_i = a$.

Recurse on that node:

- 4.1) choose $X_i = a$, a particular branch
- 4.2) calculate $H(Y|X_i=a)$.
- 4.3) calculate $H(Y|X_1, X_i=a), \dots, H(Y|X_{i-1}, X_i=a), H(Y|X_{i+1}, X_i=a), \dots, H(Y|X_d, X_i=a)$
- 4.4) calculate info gain:
 $IG(Y, X_1, X_i=a) \dots IG(Y, X_d, X_i=a)$.
- 4.5) choose attribute w/ LARGEST IG as DECISION NODE.
- 4.6) choose $X_i = b$, the second branch
!

Recurse on the next layer, repeat above ...

+ GOOD TREE:

- + "Goldilocks size" → too small won't handle subtle distinctions in data
→ too big will be efficient and overfit.
- + find smallest tree fitting hypothesis
- + INDUCTIVE BIASE: small trees, informative nodes near root

+ ISSUES:

- + less data at lower levels (exponentially)
- + ↑ size tree = overfits
- + Above algo greedy → NOT OPTIMAL b/c NP-Hard Problem.
- + can REGULARIZE = prune, min gain for split, min node size

+ PROS:

- + Good when attrs are boolean, depends critically on few attributes

+ CONS:

- + Bad on parity and most functions
- + Not good on continuous attributes

CLASS CLASSIFICATION

VR: R1\&R2 region ambiguity, OVO: no-R region ambiguity
K-CLASS DISCRIMINANT:

$$y_u(x) = w_u^T x + w_{u,0} \quad (K \text{ unique } w's)$$

Assign if $y_u(x) \geq y_j(x)$

$$(w_u - w_j)^T x + (w_{u,0} - w_{j,0}) = 0 \leftarrow \text{line}$$

+ Convex decision regions

+ SOFTMAX:

$$p(c_u | x) = y_u(x) = \frac{\exp(z_u)}{\sum_j \exp(z_j)}$$

+ ONE-HOT ENCODING: $t = [0, 1, 0, 0]^T \rightarrow u = 2, t = [t_1, t_2, \dots, t_K]$

+ MULTICLASS LOGISTIC REGRESSION

+ We want to MAXIMIZE THE LIKELIHOOD of the correct class classification.

+ Let T be a matrix of N dims, row i is the t^i

+ Then: WT Maximize:

$$p(T | X; w_1, \dots, w_K) = \prod_{i=1}^N \prod_{u=1}^K p(c_u | x_i)^{t_u^i} = \prod_{i=1}^N \prod_{u=1}^K y_u^i(x_i)^{t_u^i} \quad (\text{IID}).$$

Recall neg. log likelihood:

$$-\log p(T | X; w_1, \dots, w_K) = -\log \left[\prod_{i=1}^N \prod_{u=1}^K y_u^i(x_i)^{t_u^i} \right]$$

$$L(w_1, \dots, w_K) = - \sum_{i=1}^N \sum_{u=1}^K t_u^i \log(y_u^i(x_i)) \leftarrow \text{CROSS-ENTROPY ERROR}$$

+ Obtaining OPTIMIZED WEIGHTS

+ Do GRAD DESCENT. Need derivatives:

$$\frac{\partial y_j^i}{\partial z_u^i} = \frac{\partial}{\partial z_u^i} \left(\frac{\exp(z_j^i)}{\sum_l \exp(z_l^i)} \right) = \underbrace{\delta(u=j)}_{\substack{l=i \\ u=j, \\ 0 \text{ else}}} y_j^i - y_j^i y_u^i$$

$$\frac{\partial}{\partial z_u^i} \left(\frac{\exp(z_j^i)}{\sum_l \exp(z_l^i)} \right) = \underbrace{\frac{\partial L}{\partial z_u^i}}_{\substack{u=j \\ 0 \text{ else}}} = y_u^i - t_u^i$$

$$\frac{\partial L}{\partial w_{u,g}} = \sum_{i=1}^N \frac{\partial L}{\partial z_u^i} \cdot \frac{\partial z_u^i}{\partial w_{u,g}} = \sum_{i=1}^N (y_u^i - t_u^i) x_g^i$$

Note: $\frac{\partial L}{\partial w_{u,g}} = \sum_{i=1}^N \frac{\partial L}{\partial z_u^i} \cdot \frac{\partial z_u^i}{\partial w_{u,g}} = \sum_{i=1}^N \left(\sum_{j=1}^K \frac{\partial L}{\partial y_j^i} \cdot \frac{\partial y_j^i}{\partial z_u^i} \right) \cdot \frac{\partial z_u^i}{\partial w_{u,g}}$

+ K-NNs, DCs can directly handle multiclass!

GENERATIVE MODELS

- + Model $p(x, y)$ given $p(y)$, and compute $p(x|y)$. KEY IDEA: The probability of class y ! "HOW DOES DATA LOOK LIKE FOR A CLASS?"
- + DISC: How do I separate the classes?
- + GEN: What does each class "look" like?
- + Post = $\frac{\text{class likelihood} \times \text{prior}}{\text{evidence}}$
- + GOAL: $p(c|x) = \frac{p(x|c)p(c)}{p(x)}$
- + SIMPLE BERNoulli: Binary $p(x)$ class situation:
 - + $p(x|c) = \text{Ber}(\theta_c) = \theta_c^x (1-\theta_c)^{1-x}$
 - + FIT DISTRIBUTION TO DATA: Let D_c be data of class c .
 - + Want to MAXIMIZE $p(D_c|c) \rightarrow$ prob. of data of class c , given class c .
 - + $p(D_c|c) = \prod_{i=1}^N p(x_i|c) = \prod_{i=1}^N \theta_c^{x_i} (1-\theta_c)^{1-x_i} = \theta_c^{N_c} (1-\theta_c)^{N-N_c}$
 - (note $x_i \in \{0, 1\}$ here)
 - + MINIMIZE $-\log$: $-\log p(D_c|c) = -\log (\theta_c^{N_c} (1-\theta_c)^{N-N_c})$
 $= -N_c \log(\theta_c) - (N-N_c) \log(1-\theta_c)$
 - + DERIVATIVE FOR SGD:

$$\frac{\partial \log \theta_c}{\partial \theta_c} = -\frac{N_c}{\theta_c} + \frac{N-N_c}{1-\theta_c} = 0 \rightarrow \boxed{\theta_c = \frac{N_c}{N}}$$

- + What if $N_c = 0$? (Rare Word Case)

+ BETA-BINOMIAL: Need a PRIOR over θ .

$$\text{Beta}(\theta|a, b) = \theta^{a-1} (1-\theta)^{b-1}$$

$$+ p(\theta_c|D_c) \propto p(D_c|\theta_c) P(\theta_c) \rightarrow \theta_c^{N_c} (1-\theta_c)^{N-N_c} \cdot \theta_c^{a-1} (1-\theta_c)^{b-1}$$

$$= \theta^{N_c+a-1} (1-\theta_c)^{N-N_c+b-1}$$

$$+ \text{MAP estimation: } \hat{\theta}_c, \text{ map: } \frac{N_c+a-1}{N+a+b-2}$$

NAIVE BAYES

- + Before: IID assumption: $p(y|x; w) = p(y_1|x_1; w) \dots p(y_d|x_d; w)$
- + Now: $p(x, y; \theta_c) = p(x_1|x_2 \dots x_d; \theta_c) p(x_2|x_3 \dots x_d, y; \theta_c) \dots p(x_{d-1}|x_d, y; \theta_c) p(y; \theta_c)$
- + NEED NAIVE BAYES ASSUMPTION: Dimensions d are INDEPENDENT given class y .

$$p(x|y=c; \theta_c) = \prod_{i=1}^d p(x_i|y=c; \theta_{ic}) = \prod_{i=1}^d \theta_{ic}^{x_i} (1-\theta_{ic})^{1-x_i}$$

+ First, log it:

$$\begin{aligned} \log p(x|y=c; \theta_c) &= \sum_{i=1}^d x_i \log \theta_{ic} + \sum_{i=1}^d (1-x_i) \log (1-\theta_{ic}) \\ &= \sum_{i=1}^d x_i \log \theta_{ic} + \sum_{i=1}^d y_i \log (1-\theta_{ic}) + \sum_{i=1}^d \log (1-\theta_{ic}) \end{aligned}$$

$$p(x|y=c; \theta_c) = \exp(\log p)$$

$$= \sum_{i=1}^d x_i (\log \theta_{ic} / 1-\theta_{ic}) + \sum_{i=1}^d \log (1-\theta_{ic})$$

"define $w_{ci} = \log(\theta_{ic}/1-\theta_{ic})$, $w_{0c} = \sum_{i=1}^d \log(1-\theta_{ic})$
 $P(x|y=c; \theta_c) = \exp(w_0^T x + w_{0c})$

HOW TO CLASSIFY:

$P(y=c|x) \propto P(x|y=c; \theta_c) P(y=c|\theta_c)$.

Thus,

$P(y=c|x) = \exp[w_0^T x + w_{0c} + \log(P(y=c|\theta_c))]$

LINEAR CLASSIFIER!

Really easy to train. NO DERIVATIVES. $|w_{ci}|$ is just COUNTING. $O(d)$.

+ $p(y=c|\theta_c) = \frac{\# \text{ class } C}{N}$ classifier by max_c $P(y=c|x)$

- + PROS:
 - + Fast train, test
 - + Less overfit, can be better than logistic on small data
 - + Easy to add/remove classes
 - + Can handle partial data

- + CONS: + Naive IID doesn't always hold!

GAUSSIAN MIXTURE MODELS + EM

- + GOOD CLUSTERING \rightarrow assume data generated by generative model!
 \rightarrow adjust MODEL PARAMS to maximize probability that the model produces the data we observe.

- + Usually: $p(x, z) = p(x|z)p(z)$, but we don't know $z = ?$

so, MIXTURE MODEL:

$$p(x) = \sum_z p(x, z) = \sum_z p(x|z)p(z).$$

- + GMM: K classes, $p(z=u) = \pi_u$, $p(x|z=u) = N(x|\mu_u, \Sigma_u)$

then, $p(x) = \sum_{u=1}^K N(x|\mu_u, \Sigma_u) \pi_u$ and $\sum_u \pi_u = 1$

+ UNIVERSAL DENSITY APPROXIMATOR!

- + MAXIMUM LIKELIHOOD:

$$\log p(x|\pi, \mu, \Sigma) = \log \left(\prod_{i=1}^N p(x_i) \right) = \sum_{i=1}^N \log \left(\sum_{u=1}^K N(x_i|\mu_u, \Sigma_u) \pi_u \right).$$

+ SINGULARITIES: \uparrow likelihood when Gaussian explains single pt.

+ NONCONVEX

+ Soln. INVARIANT to permutes

- + LATENT VARIABLES: z^i for each x^i .

Then, $\log p(x|\pi, \mu, \Sigma) = \sum_{i=1}^N \log \left(\sum_{u=1}^K p(x_i|z^i; \mu_u, \Sigma_u) \pi_u \right)$

We don't know z^i for every x^i .

(else: $\mu_u = \frac{1}{S(z^i, u)} \sum_{i=1}^N \delta(z^i, u)x^i$ (mean of all u -x's)

$$\Sigma_u = \frac{1}{S(z^i, u)} \sum_{i=1}^N \delta(z^i, u) (x^i - \mu_u)(x^i - \mu_u)^T$$
 (covar.)

$$\pi_u = \frac{1}{N} \sum_{i=1}^N \delta(z^i, u). \quad (\text{all } u\text{'s over-}\Delta)$$

+ IDEA: The EM algorithm to OPTIMIZE PARAMETERS

1) E: compute POST. PROB OVER z GIVEN CURRENT MODEL.
PARAMS. (How much do we think each Gaussian
gens each datapt.)

M: Assume data was generated this way.
change PARAMS of each Gaussian to MAXIMIZE PROB
that it would gen. the data its resp. for

E-STEP:

1) Compute post. prob over z .

$$y_u^i = P(z_i=k | x_i; \pi, \mu, \Sigma),$$

"which Gaussian gen'd each datapt?" AREN'T SURE \rightarrow DISTRIBUTION OVER ALL POSSIBLE GAUSSIANS (u).

M-STEP:

1) Each Gaussian now has some posterior prob for each data point,
which is y_u^i

2) FIT each GAUSSIAN to the WEIGHTED DPTS.

3) Derive closed-form UPDATES OF ALL PARAM

WHERE DOES EM COME FROM?

likelihood.

+ Let $\Theta = \mu, \pi, \Sigma$.

$$+ l(x, \Theta) = \sum_i \log \left(\underbrace{\sum_u P(x_i, z_i=k; \Theta)}_{q_u} \right) \leftarrow \text{hard!}$$

+ A TRICK IN ML: INTRODUCE NEW DISTRIBUTION q:

$$l(x, \Theta) = \sum_i \log \left(\sum_u q_u \frac{P(x_i, z_i=k; \Theta)}{q_u} \right)$$

+ JENSEN'S INEQUALITY: (concave fun like log).
 $f(\mathbb{E}[x]) \geq \mathbb{E}[f(x)]$

$$f\left(\sum_i p_i x_i\right) \geq \sum_i p_i f(x_i)$$

maximize this =
lower bd = increase likelihood!

+ Apply JENSEN to the above:

$$l(x, \Theta) = \sum_i \log \left(\sum_u q_u \frac{P(x_i, z_i=k; \Theta)}{q_u} \right) \geq \sum_i \sum_u q_u \log \left(\frac{P(x_i, z_i=k; \Theta)}{q_u} \right)$$

+ Define $q_u = p(z_i=k | x_i; \Theta^{old})$.

$$Q(\Theta) = \sum_i \sum_u p(z_i=k | x_i; \Theta^{old}) \log (P(x_i, z_i=k; \Theta)),$$

$$= \mathbb{E}_{P(z_i=k | x_i; \Theta^{old})} [\log P(x_i, z_i=k; \Theta)] \quad \leftarrow$$

+ CONCEPTUALLY: we don't know z_i , so we AVERAGE them given the CURRENT MODEL Θ^{old} .

E-STEP:

Calculating POST PROB RESPONSIBILITIES:

$$\begin{aligned}
 Y_u = p(z=k|x) &= \frac{p(x|z=k)p(z=k)}{p(x)} \\
 &= \frac{p(x|z=k)p(z=k)}{\sum_{j=1}^k p(x|z=j)p(z=j)} \quad \left. \right\} \text{Bases} \\
 &= \frac{N(x|\mu_u, \Sigma_u)\pi_u}{\sum_{j=1}^k N(x|\mu_j, \Sigma_j)\pi_j} \quad \left. \right\} \text{Gaussian MM.}
 \end{aligned}$$

+ Now we have $\gamma_u^i = p(z^i=k|x)$.

+ Calculating EXPECTED LIKELIHOOD:

$$\begin{aligned}
 E_{p(z^i|x^i)} \left[\sum_i \log(p(x^i|z^i; \theta)) \right] &= \\
 &= E_{p(z^i|x^i)} \left[\sum_i \log(p(x^i|z^i; \theta)p(z^i; \theta)) \right] \\
 &= \sum_i \sum_u \gamma_u^i \log(p(x^i|z^i=k; \theta)) + \gamma_u^i \log(p(z^i=k; \theta)) \\
 &= \sum_i \sum_u \gamma_u^i \left(\log(N(x^i; \mu_u, \Sigma_u)) + \log(\pi_u) \right) \\
 &= \sum_u \sum_i \gamma_u^i \log(N(x^i; \mu_u, \Sigma_u)) + \sum_u \sum_i \gamma_u^i \log(\pi_u).
 \end{aligned}$$

GMM M-STEP

- + Need to optimize $\sum_u \sum_i \gamma_u^i \log(N(x^i; \mu_u, \Sigma_u)) + \sum_u \sum_i \gamma_u^i \log(\pi_u)$
- + Solving for μ_u, Σ_u is LINE FITTING w/ SEP GAUSSNS weighted γ_u^i .
- + Similar to generate derivation.

ALGO:

- 1) Init π_u, Σ_u, μ_u
- 2) ITERATE TILL CONVG:
 - 3) E-STEP: calculate post probs $\gamma_u^i = p(z^i=k|x)$.
 - 4) M-STEP: Re-estimate params w/ new responsibilities

$$\mu_u = \frac{1}{N_u} \sum_{i=1}^N \gamma_u^i x^i$$

$$\Sigma_u = \frac{1}{N_u} \sum_{i=1}^N \gamma_u^i (x^i - \mu_u)(x^i - \mu_u)^T$$

$$\pi_u = \frac{N_u}{N}, \quad N_u = \sum_i \gamma_u^i$$
- 5) Eval log-likelihood check for convergence.

K-MEANS

- + Why does it converge?
 - + Sum of all distances between pts and means J always REDUCES, and always start pos, so will convg. to 0
 - + CONVERGENCE → assignments do not change
- + Local Minima?
 - + J is NONCONVEX, can get stuck in local minima
 - + FIX:
 - multiple random starts
 - split big cluster in 2, merge two nearby clusters
- + K-Means++:
 - + Randomly choose first mean/center
 - + If x^i s.t. $d_i = \text{dist. closest to center}$
 - + Pick new center to be at x^i w/ prob $d_i^{(in)^2}$
 - + Repeat till k centers.
- + Soft K-Means: ALGO:
 - 1) Set μ_u to random
 - 2) ITERATE till conv:
 - 3) $\forall x^i$, check all μ_u :
 - 4) Give $r_{u^i} = \frac{\exp(-\beta d(\mu_u, x^i))}{\sum \exp(-\beta d(\mu_u, x^j))}$ } SOFTMAX ASSIGNMENT
 - 5) Update model params
$$\mu_u = \frac{\sum r_{u^i} x^i}{\sum r_{u^i}}$$
 (like hard k-Means)

LE METHODS

S-VARIANCE DECOMP:

- + regression w/ L_2 loss, $h^*(x) = \mathbb{E}[t|x]$.
- + ML prediction: $y(x; D)$

+ Breaking $\mathbb{E}_{D,x,t}[(t - y(x; D))^2]$ into core components:

$$\begin{aligned}\mathbb{E}_{D,x,t}[(t - y(x; D))^2] &= \mathbb{E}_{D,x,t}[(t - h^*(x) + h^*(x) - y(x; D))^2] \\ &= \mathbb{E}_t[(t - h^*(x) + h^*(x) - y(x; D))^2 | x, D] \\ &= \mathbb{E}_t[(t - h^*(x))^2 | x, D] + \mathbb{E}_t[(h^*(x) - y(x; D))^2 | x, D] \\ &\quad + \underbrace{2\mathbb{E}_t[(t - h^*(x))(h^*(x) - y(x; D))]}_{2(\mathbb{E}_t[t|x] - h^*(x))(h^*(x) - y(x; D))} \\ &= 2(h^*(x) - \cancel{h^*(x)})(h^*(x) - y(x; D))\end{aligned}$$

Thus:

$$\mathbb{E}_{D,x,t}[(t - y(x; D))^2] = \underbrace{\mathbb{E}_t[(t - h^*(x))^2 | x, D]}_{\text{error}} + \underbrace{\mathbb{E}_t[(h^*(x) - y(x; D))^2 | x, D]}_{\text{noise}}$$

$$\begin{aligned}\text{Breakdown: } \mathbb{E}_D[(h^*(x) - y(x; D))^2 | x] &= \mathbb{E}_D[(h^*(x) - \mathbb{E}_D[y(x; D)] + \mathbb{E}_D[y(x; D)] - y(x; D))^2 | x] \\ &= \mathbb{E}_D[(h^*(x) - \mathbb{E}_D[y(x; D)])^2 | x] + \mathbb{E}_D[(\mathbb{E}_D[y(x; D)] - y(x; D))^2 | x] \\ &\quad + \underbrace{2\mathbb{E}_D[(h^*(x) - \mathbb{E}_D[y(x; D)])(\mathbb{E}_D[y(x; D)] - y(x; D))]}_{2(h^*(x) - \mathbb{E}_D[y(x; D)])(\mathbb{E}_D[y(x; D)] - \cancel{\mathbb{E}_D[y(x; D)]})} | x\end{aligned}$$

Thus,

$$\mathbb{E}_{D,x,t}[(t - y(x; D))^2] = \underbrace{\mathbb{E}_t[(t - h^*(x))^2]}_{\text{error}} + \underbrace{\mathbb{E}_x[(h^*(x) - \mathbb{E}_D[y(x; D)])^2]}_{\text{noise}} + \underbrace{\mathbb{E}_{D,x}[(y(x; D) - \mathbb{E}_D[y(x; D)])^2]}_{(\text{bias})^2}$$

- + HIGH VARIANCE (OVERFITS) \rightsquigarrow LOW BIAS
- + LOW VARIANCE \rightsquigarrow HIGH BIAS

+ BAGGING:

- + TO REDUCE VARIANCE \approx OVERFITTING
- + BOOTSTRAPPING: sample w/ replacement from data to get variance estimate

+ ALGO:

INPUT: D , ML algo A , # bags N .

1) ITER $i: 1 \rightarrow N$:

- 2) Get $D_i \subseteq D$ by SAMPLING w/ REPLACEMENT
- 3) $f_i = A(D_i)$.

4) Return f_1, \dots, f_N

5) PREDICTION: return $\frac{1}{N} \sum_i f_i(x)$ or majority vote.

+ RANDOM FOREST

- + Random subset of features on each split \rightarrow reduce correlation.
- + D-Trees: easily overfit, fast inference, good on relational data.
- + OOB EST: cheaply estimate test loss \rightarrow can reduce validation need
 - + PREDICT each TRAINING EXAMPLE using ALL TREES that DID NOT CONTAIN THIS EXAMPLE in training data.
 - + Slightly worse loss than real test error \rightarrow upper bound.
- + Reduces overfit by avg. predictions
- + DOES NOT REDUCE BIAS. STILL CLASSIFIER CORRELATION.
- + RF soln: MORE RANDOMNESS

BOOSTING

Reduces BIAS.

- + Idea: Take M WEAK CLASSIFIERS that does ϵ better than chance and BOOST IT to get LOW TRAINING ERROR.

ADABOOST ALGO:

INPUT: $\{x_i, t_i\}_{i=1}^N$, WeakLearn learning procedure.

INIT: D is an $N \times M$ matrix where each column D_m is a set of weights for the N examples. Initialize $D_i^j(x) = 1/N$ for all i to N , j 's.

$$\begin{bmatrix} 1/N & \dots & 0 \\ 1/N & \dots & 0 \\ \vdots & \ddots & \vdots \\ 1/N & \dots & 0 \end{bmatrix}$$

1) ITER: $m = 1 \rightarrow M$

2) $h_m(x) = \text{WeakLearn}(\{x\}, t, w)$. Get weak classifier

3) Minimize:

$$J_m = D_m^1[h_m(x^1) \neq t^1] + D_m^2[h_m(x^2) \neq t^2] + \dots + D_m^N[h_m(x^N) \neq t^N].$$

(Idea: J_m is the "score" of misclassifications by the weak classifier weighted by D_m).

4) Compute WEIGHTED ERROR RATE:

$$E_m = \frac{J_m}{\sum_i D_m^i} \quad \left\{ \text{always } < 1/2 \right.$$

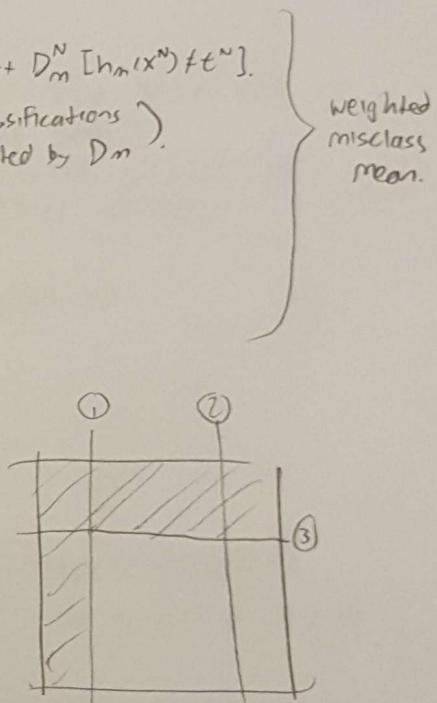
5) Compute CLASSIFIER COEFF: $\alpha_m = \frac{1}{2} \log \frac{1-E_m}{E_m}$.

6) UPDATE DATA WEIGHTS:

$$D_{m+1}^i = D_m^i \exp(-t^i \alpha_m h_m(x^i)).$$

7) FINAL:

$$H(x) = \text{sign}\left(\sum_{m=1}^M \alpha_m h_m(x)\right) = \text{sign}(F(x))$$



+ BOOSTING \rightarrow greedy optimization of exponential loss