

### Problem 1 - Probability

(a) Given: • sample  $\vec{x} = (x_1, \dots, x_n)$   
• Bernoulli distribution  $\rightarrow$  Find: maximum likelihood estimator of  $p$

Let's say probability of  $(y=1|\vec{x}) \equiv \phi$

$\therefore P(y=1|\vec{x}) = \phi^k (1-\phi)^{(1-k)}$  where  $k \in \{0, 1\}$    
 belongs to range

$\therefore$  likelihood  $L(\phi) = \prod_{i=1}^n [\phi^{x_i} (1-\phi)^{(1-x_i)}]$

$\therefore$  log likelihood  $\ell(\phi) = \log \phi \sum_{i=1}^n x_i + \log(1-\phi) \sum_{i=1}^n (1-x_i)$

$\rightarrow \frac{\partial \ell(\phi)}{\partial \phi} = \frac{\sum x_i}{\phi} - \frac{\sum (1-x_i)}{1-\phi} \stackrel{\text{set}}{=} \text{zero}$

$\xrightarrow{\text{multiply by } \phi(1-\phi)} (\sum x_i) - \hat{\phi}(\sum x_i) = \hat{\phi} \sum (1-x_i) \rightarrow \hat{\phi} n = \sum x_i$

$\therefore \hat{\phi} = \frac{1}{n} \sum_{i=1}^n x_i = \bar{x}$

since  $\frac{\partial^2 \ell(p)}{\partial p^2}$  is negative, we confirm this is a maximization

(b) Now suppose the distribution is uniform

Find: MLE:

$$P(y_i | \vec{x}; a, b) = \begin{cases} 0 & \text{for } x < a \\ 1/(b-a) & \text{for } a \leq x \leq b \\ 0 & \text{for } x > b \end{cases}$$

$$\therefore \text{likelihood } L(a, b | \vec{x}) = \prod_{i=1}^n \frac{1}{b-a}$$

$$\therefore \ell(\cdot) = \sum_{i=1}^n \log\left(\frac{1}{b-a}\right) = \sum_{i=1}^n \left( \log(1) - \log(b-a) \right)$$

$$\xrightarrow{\text{partial to } a} \frac{\partial \ell(\dots)}{\partial a} = \sum_{i=1}^n \left( 0 - \frac{-1}{b-a} \right) \therefore \hat{a} = \min(\vec{x})$$

$$\xrightarrow{\text{partial to } b} \frac{\partial \ell(\dots)}{\partial b} = \sum_{i=1}^n \left( 0 - \frac{1}{b-a} \right) \therefore \hat{b} = \max(\vec{x})$$

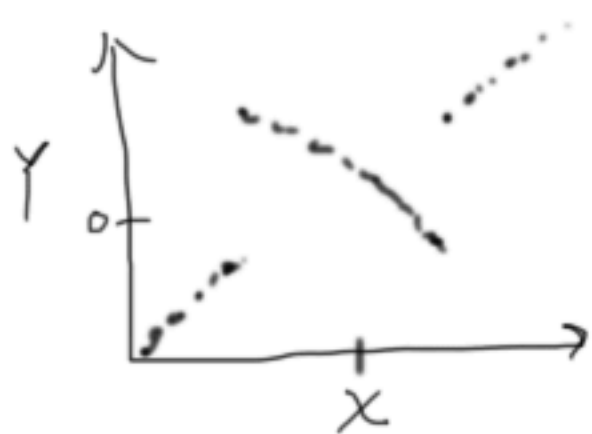
$\therefore$  we want both to approach zero  
(Both 2<sup>nd</sup> derivatives are negative) ✓

(c) Given: • X, Y random  
• both marginally Gaussian  
• correlation = 0

Find: Are they independent?

given  $\nRightarrow$  independence

Counter example: X maps to  $f(z) = z$   
Y maps to  $g(z) = \begin{cases} f(z) & \text{for } |z| \leq c \\ -f(z) & \text{for } |z| > c \end{cases}$



Y depends on X, yet the raw correlation equals zero by the average, by some c (not all).

X has a gaussian mean and variance, as does Y, by its cor.

This all considered, given  $\nRightarrow$   $\neg$  independent.  
 $\therefore$  need more info.

Intermediate value theorem  
Continuous function  
value between  
that makes  
sense

## Problem 2 - Poisson GLM's

- Given:
- sample pairs  $\vec{xy} = \{(x_1, y_1), \dots, (x_n, y_n)\}$
  - $x_i \in \mathbb{R}^p$  ( $p$  dimensions)
  - $y_i \in \mathbb{N} \{0, 1, 2, 3, \dots\}$  for  $i=1, \dots, n$  (countable)
  - linear coefficients are  $\vec{\theta}$ , also then in  $\mathbb{R}^p$

(a) Find: log likelihood  $\ell$  aka  $LL$

For GLM form  $\lambda$  in Poisson is  $e^{\vec{\theta}^T \vec{x}}$  (via  $\eta$  in exponential family form)

$$L(\vec{\theta} | \vec{xy}) = \prod_{i=1}^n \left[ \frac{e^{y_i \vec{\theta}^T x_i} e^{-e^{\vec{\theta}^T x_i}}}{y_i!} \right] \because p(y_i | x_i) = \frac{\lambda^{y_i}}{y_i!} e^{-\lambda}$$

$$\xrightarrow{\text{take log}} LL = \sum_{i=1}^n \log(L(\vec{\theta} | \vec{xy})) = \sum_{i=1}^n (\log(e^{y_i \vec{\theta}^T x_i} e^{-e^{\vec{\theta}^T x_i}}) - \log(y_i!)) = \sum_{i=1}^n (y_i \vec{\theta}^T x_i - e^{\vec{\theta}^T x_i} - \log(y_i!))$$

$$\nabla LL(\vec{\theta} | \vec{xy}) \stackrel{\text{set}}{=} \text{zero} \rightarrow \frac{\partial LL(\vec{\theta} | \vec{xy})}{\partial \vec{\theta}} = \sum y_i x_i - x_i e^{\vec{\theta}^T x_i} \text{ (no closed form soln.)}$$

(b) Assuming a MLE  $\hat{\theta}$  is solved, predicting a  $y$  given  $x^*$  follows the form:

The probability of a  $y$  given  $x^*$  w/ parameters  $\hat{\theta}$

$$p(y | x^*; \hat{\theta}) = \frac{\lambda^y}{y!} e^{-\lambda} = \frac{e^{y \hat{\theta}^T x^*} e^{-e^{\hat{\theta}^T x^*}}}{y!}$$

(c)  $x^*$  contains a dimension not in  $\mathbb{R}^p$ , returning  $x^* \cdot \hat{\theta}^T = 0$ .

Predicted  
 $\eta = \log \lambda$

$\lambda = \theta^T x$

$\hat{\theta}^T x^*$

To find the most likely  $y$ , iterate from  $y=0$  to  $n-1$  and return the max from that search (or set  $\nabla \rightarrow 0$ )

original prediction of  $y \Rightarrow \frac{e^{y(0)} e^{(0)}}{y!} = \frac{1}{y!}$ , which is unintended.

L2 ridge regression adds  $\alpha \|\theta\|^2$   
to differentiate from Poisson?

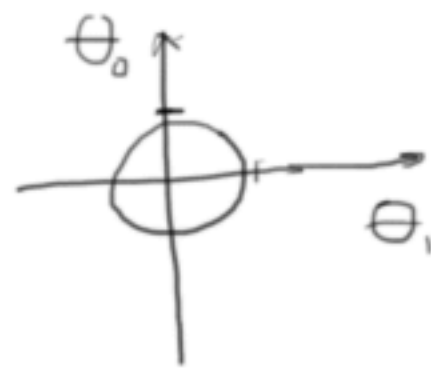
If we are predicting a maximum likelihood  $\hat{\theta}$  first (with  $\vec{x}_y'$  being  $\vec{x}_y$  with  $x^*$ )

$$\text{loss func } \mathcal{L}_2(\theta | \vec{x}_y') = \sum_{i=1}^n (y_i \vec{\theta}^T x_i - e^{\vec{\theta}^T x_i} - \log(y_i)) + \alpha \|\theta\|^2$$

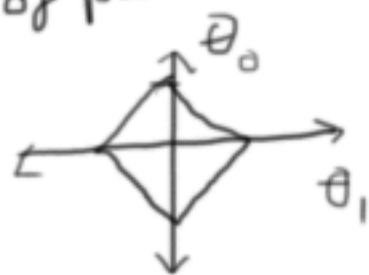
$$\nabla \mathcal{L} = \frac{\partial \mathcal{L}}{\partial \theta} = \sum_{i=1}^n ((y_i x_i - x_i e^{\vec{\theta}^T x_i} - 0)) + 2\alpha \theta = \begin{bmatrix} \frac{\partial \mathcal{L}}{\partial \theta_0} \\ \vdots \\ \frac{\partial \mathcal{L}}{\partial \theta_j} \end{bmatrix}$$

For  $x_i = x^*$ ,  $\mathcal{L}_{x^*} = y_i x^* - x^* e^{(\cdot)} = y_i x^* - x^*$ , also unintended.

Additionally, there are not enough data points with the orthogonal dimension to estimate a sensible weight in  $\vec{\theta}$ , but the weight will approach to 0, never to 0 itself, and affect all other predictions.



(d) Since we seek to bring the weight of an irrelevant/untenable  $\theta_j$  to absolute 0, L1 lasso regression offers this by penalizing parameters that do not significantly help the loss function. It trades off inclusion for relevance, by requiring a selection of parameters.



### Problem 3 - Very Random Forest

Given;  $\mathcal{D} = \{(x_1, y_1), \dots, (x_n, y_n)\}$

- $x_i \in \mathbb{R}^p$
- $y_i \in \mathbb{R}$
- sample one feature
- sample  $m < n$   $x_j$ 's ( $\tilde{x}$ )

(a) Find: Data distributions where this "very" random forest has no bias.

Bias is expected difference between:

$$E(y|x) \quad \text{vs} \quad E_{\text{rf}}[f(x); \mathcal{D}]$$

squared.

$\therefore$  the single feature trees can account for functions with terms that involve a single feature each in the true distribution, because the ensemble averages the single feature trees together. Assuming  $\infty$  trees.

$$\text{e.g. } f(\vec{x}) = x_1 + x_2 + \dots + x_n \quad \text{trees}$$

$$\text{or } f(\vec{x}) = x_1^{k_1} + x_2^{k_2} + \dots + x_n^{k_n} \quad \text{where } k_i \in \mathbb{R}$$



(b) With an infinite amount of trees, results in both regular RF and VRF can theoretically reach the truth of  $\mathcal{D}$ 's distribution.

|     | BIAS | VARIANCE |
|-----|------|----------|
| RF  | none | higher   |
| VRF | none | lower    |

Further discussion on bias and variance.

$$\text{Variance} \equiv \mathbb{E} [(\text{longterm-performance} - \text{expected test performance})^2]$$

$$= \mathbb{E}_{g \sim p} [(\mathbb{E}_{p \sim p} [f(x; \mathcal{P})] - f(x; \mathcal{P}))^2]$$

With an infinite amount of trees, the long term performance should closely capture the true distribution described. Test performance ensembled from single feature trees will be accurate, as trees  $\rightarrow \infty$ , of course if data captures (will be noise). However, traditional RF trees are less stable in terms of this distribution, since they can introduce small interfeature dependencies which do not actually exist. But with trees approaching infinity, variance too should approach zero (again, only in distributions where vRF has 0 bias.)

We can say vRF approaches 0 variance faster than RF.  
In both cases, forest will overfit if not given new data, thus leading to variance.  
In a diagram:



$f^* \equiv$  true function

$M_{vRF} \equiv$  model of very Random Forest  $\propto$

$M_{RF} \equiv$  model of traditional Random Forest  $\propto$

$M_{RF}$  more robust but can have more overall variance  
 $M_{vRF}$  = less robust and more likely to have small variance given that answer is in model space.

Trivially, uniform distributions also have 0 bias, and variance for both models, given  $x^* \in \text{domain}$ .

#### Problem 4 - Alternative Losses

- Given:
- $\vec{x}, \vec{y} \{ (x_1, y_1), \dots, (x_n, y_n) \}$
  - $x_i \in \mathbb{R} \quad y_i \in \{0, 1\}, i=1, \dots, n$
  - $P(y_i = 1 | x_i) = \sigma(\beta x_i) = \frac{1}{1 + e^{-\beta x_i}}$
  - $P(y_i = 0 | x_i) = 1 - \sigma(\beta x_i)$

(a) Find: derivative of squared loss wrt  $\beta$  at  $\beta = 3$

$$\frac{\partial \sigma(\beta x_{n+1})}{\partial \beta} = \left[ y_{n+1} - \sigma(\beta x_{n+1}) \right]^2 \quad \begin{cases} x_{n+1} = 100 \\ y_{n+1} = 0 \end{cases}$$

$$= 2 \left( y_{n+1} - \frac{1}{1 + e^{-\beta x_{n+1}}} \right) \frac{\partial}{\partial \beta} \left[ y_{n+1} - \frac{1}{1 + e^{-\beta x_{n+1}}} \right]$$

$$= \left( 2 x_{n+1} e^{\beta x_{n+1}} \left( (y_{n+1} - 1) e^{\beta x_{n+1} + y_{n+1}} \right) \right) / (e^{\beta x_{n+1}} + 1)^3$$

$$= \frac{2(100) e^{3(100)} (-e^{300})}{(e^{300} + 1)^3} = 1.029 \times 10^{-128} \text{ for } P(y=1)$$

negative for  $P(y=0)$



$$(b) \mathcal{L}(\beta) = \prod_{i=1}^n p(Y=y_i | X=x_i) = \prod_{i=1}^n \sigma(\beta x_i)^{y_i} \cdot (1 - \sigma(\beta x_i))^{(1-y_i)}$$

$$\therefore \underset{\text{log likelihood}}{\ell(\beta)} = \sum_{i=1}^n y_i \log \sigma(\beta x_i) + (1-y_i) \log (1 - \sigma(\beta x_i))$$

$$\underset{\text{loss}}{\mathcal{L}} \equiv -\ell(\beta) \therefore \frac{\partial \mathcal{L}(\beta)}{\partial \beta_j} \Rightarrow \text{chain rule} \quad \text{where } \sigma_{\beta x} \equiv \sigma(\beta x) \quad \frac{\partial \mathcal{L}(\beta)}{\partial \sigma_{\beta x}} \cdot \frac{\partial \sigma_{\beta x}}{\partial \beta_j}$$

$$\text{echoing the logit function} \quad \log\left(\frac{p}{1-p}\right) \quad \Rightarrow \text{chain rule} \quad \text{where } \gamma \equiv \beta x \quad \frac{\partial \mathcal{L}(\beta)}{\partial \sigma_{\beta x}} \cdot \frac{\partial \sigma_{\beta x}}{\partial \gamma} \cdot \frac{\partial \gamma}{\partial \beta_j}$$

$$\Rightarrow \text{Substitution with derived results} \quad -\left(\frac{y}{\sigma_{\beta x}} - \frac{1-y}{1-\sigma_{\beta x}}\right) \cdot \sigma_{\beta x}(1-\sigma_{\beta x}) \cdot x_j$$

$$\Rightarrow \text{expanding} \quad -(y(1-p) - p(1-y)) \cdot x_j$$

$$\Rightarrow -(y-p) x_j = -(y - \sigma(\beta x)) x_j$$

$$\Rightarrow \underset{\substack{x=100 \\ y=0 \\ \beta=3}}{-\left(0 - \frac{1}{1+e^{3(100)}}\right)(100)} \approx \underset{\text{for } \sigma(\beta x)}{-(0-1)(100)} = 100$$

(c) These loss functions imply drastically different characteristics of the parameter  $\beta$ . This is because different distributions call for different concepts of normalization. For some distributions, such as Bernoulli, least squares does not make the most sense because all  $y$  points will be a set distance from the line (0 or 1).

Bernoulli uses logistic regression as seen in this problem.

This is why least squares yields an arbitrary  $\sim 0$  derivative.

In this case,  $\mathcal{L} \equiv -\ell$  yields a more palpable result.  $\therefore$  do not take  $\mathcal{L}$  for granted!

