

Problemset 1

$x \in \mathbb{R}^{n \times m}, y \in \mathbb{R}^m$

$$= x^i g(1-g), g(x) = \frac{1}{1+e^{-\theta^T x}}$$

1.

a) $\nabla_{\theta}^2 J = \nabla_{\theta} \left[\frac{\partial J}{\partial g} \nabla_{\theta} g \right] = \nabla_{\theta} \left[\left(-\frac{1}{n} \sum_{i=1}^n \frac{g^i}{g^i} - \frac{1-g^i}{1-g^i} \right) \times (1-g^i) \dot{g}^i \right]$

$$= \nabla_{\theta} \left[+ \frac{x^T}{n} \sum_{i=1}^n \frac{g^i - g^i}{(1-g^i)} + (1-g^i) \dot{g}^i \right] = \frac{x^T x}{n} \sum_{i=1}^n g^i (1-g^i) \ddot{g}^i = H \ddot{g}$$

=

The Hessian is posd because $g(1-g) \in [0,1], 2g_i - 1 \geq 0$
and because $x^T x$ is posd: $\langle v^T x, x^T v \rangle = \sum_{i=1}^n \langle v, x^i \rangle \langle x^i, v \rangle \geq 0$.

b) Done.

c) Bayes rule implies $p(y=1|x) = \frac{p(x|y=1)p(y=1)}{p(x|y=0)p(y=0) + p(x|y=1)p(y=1)}$. Plugging in
the above equations yields:

$$= \frac{\exp(-1/2(x-\mu_0)^T \dots)}{\exp(-1/2(x-\mu_0)^T \dots) \cdot (1-\alpha) + \exp(-1/2(x-\mu_1)^T \dots) \cdot \alpha} = \frac{1-\alpha}{1-\alpha + \exp[-1/2(x-\mu_0)^T \dots - 1/2(x-\mu_1)^T \dots]} + 1$$

$$= \left[\exp \left(\log \frac{1-\alpha}{\alpha} - 1/2 \left(-x^T \Sigma^{-1} \mu_0 - \mu_0^T \Sigma^{-1} x + x^T \Sigma^{-1} \mu_1 + \mu_1^T \Sigma^{-1} x + \mu_0^T \Sigma^{-1} \mu_0 - \mu_1^T \Sigma^{-1} \mu_1 \right) \right) + 1 \right]$$

Define: $\Theta_0 = \log \frac{\alpha}{1-\alpha} + \frac{1}{2} (-\mu_1^T \Sigma^{-1} \mu_1 + \mu_0^T \Sigma^{-1} \mu_0)$ and $\Theta^T = (-\mu_0 + \mu_1)^T \Sigma^{-1}$

$$\Rightarrow p(y=1|x) = \left[\exp(-(\Theta_0 + \Theta^T x)) + 1 \right]^{-1}, \text{ using that } x^T \mu_K = \mu_K^T x.$$

The decision boundary has the eq $\Theta^T x + \Theta_0 = 0$, which is clearly linear in x^T !
@ threshold 0.5

Homework Set 1 Cont

1d) $\psi, \mu_0, \mu_1, \Sigma = \operatorname{argmax} l(\tilde{\delta}, \tilde{\mu}_0, \tilde{\mu}_1, \tilde{\Sigma}) \Leftrightarrow \begin{cases} p(x|y) = \frac{\exp\left[-\frac{1}{2}(x - \mu_y)^T \Sigma^{-1} (x - \mu_y)\right]}{(2\pi)^{d/2} |\Sigma|^{1/2}}, \\ p(y) = \delta^y (1-\delta)^{1-y}, \quad y \in \{0, 1\} \\ l = \sum_{i=1}^n \log p(x|y; \dots) + \log p(y; \psi). \end{cases}$

$$\begin{aligned} l &\sim -\frac{1}{2} \log |\Sigma| + \frac{1}{2} (x - \mu_{y^i})^T \Sigma^{-1} (x - \mu_{y^i}) + y^i \log \delta + (1-y^i) \log (1-\delta) \\ &= -n \log |\Sigma| - \sum_{i=1}^n (x - \mu_{y^i})^T \Sigma^{-1} (x - \mu_{y^i}) + 2y^i \log \delta + 2(1-y^i) \log (1-\delta). \end{aligned}$$

$$\frac{\partial l}{\partial \delta} = 0 \rightarrow \sum_{i=1}^n \frac{y^i}{\delta} = \sum_{i=1}^n \frac{1-y^i}{1-\delta} \Leftrightarrow (1-\delta)/n \langle y^i \rangle = (n - n \langle y^i \rangle) \delta$$

$$\Leftrightarrow \underline{\langle y^i \rangle} = \underline{\delta} \Leftrightarrow \underline{\delta} = \frac{\sum_{i=1}^n y^i}{n} = \frac{\sum_{i=1}^n \mathbf{1}\{y^i=1\}}{n}$$

$$\frac{\partial l}{\partial \mu_k}_{k \in \{0, 1\}} = 0 \rightarrow \sum_{i=1}^n \mathbf{1}\{y^i=k\} \cdot 2 \Sigma^{-1} (x^i - \mu_k) = 0 \Rightarrow \underline{\mu_k} = \frac{\sum_{i=1}^n \mathbf{1}\{y^i=k\} x^i}{\sum_{i=1}^n \mathbf{1}\{y^i=k\}}$$

For $\frac{\partial l}{\partial \Sigma}$, two important identities are necessary: $\log |\Sigma|^2 = (\Sigma^{-1})^T$

$$(Z^T \Sigma Z)^{-1} = -\frac{1}{2} \text{tr} \Sigma^{-1} Z Z^T \Sigma^{-1}$$

$(Z^T \Sigma Z)^2 = -\text{tr} \Sigma^{-1} Z Z^T \Sigma^{-1} Z Z^T \Sigma^{-1}$. Using them, we get: $\frac{\partial l}{\partial \Sigma} = -n(\Sigma^{-1})^T + \sum_{i=1}^n \text{tr} (\Sigma^{-1} (x^i - \mu_{y^i}) (x^i - \mu_{y^i})^T \Sigma^{-1})$

use that $Z^T \Sigma = \Sigma$ and that $\text{tr} y = \text{scalar} = \text{the scalar}$

$$\Leftrightarrow n \sum_{i=1}^n \Sigma^{-1} (x^i - \mu_{y^i}) (x^i - \mu_{y^i})^T \Sigma^{-1} \Leftrightarrow \underline{\Sigma} = \frac{\sum_{i=1}^n (x^i - \mu_{y^i}) (x^i - \mu_{y^i})^T}{n}$$

Done ..!

Homework set 1 cont

- 1f) Overall, the two decision boundaries look rather similar in some interrupt and \approx similar slope.
- 1g) Performance between GDA & log. reg. is practically identical for dataset 2. I'd say GDA performs ~~fairly~~ best on dataset 2, based on what I can see with my eye. This might be due to the data in dataset 2 actually coming from two normal distributions, one for each class, thus satisfying the GDA assumptions. This notion is substantiated by the GDA & lin. reg. giving so similar results, lending credence to the sigmoid assumption of the GDA hypothesis being correct. Further, to my eye, the data here does look Gaussian for both classes.
(With the eye)
- 1h) Dataset 1 looks like $X_2 \sim \exp X_1$. Taking the log of X_2 and running GDA on the resulting data gives me a much better classifier, from what my eye can judge.

2a) Done. & 2b) Done.

c) $p(t^i=0|y^i=1, x^i) = p(y^i=1|t=0) \frac{p(t^i=0)}{p(y^i=1, x^i)} = 0 \Rightarrow p(t^i=1|y^i=1, x^i) = 1$

d) $p(t^i=1, x^i) = \underbrace{p(t^i=1|y^i=1, x^i)}_{=1} p(y^i=1) + \underbrace{p(t^i=1|y^i=0, x^i)}_{=p(y^i=0|t^i=1)} \underbrace{p(y^i=0, x^i)}_{=1-\alpha}$
 $= p(y^i=1) + (1-\alpha)p(t^i=1|x^i)$
 $\Leftrightarrow \underbrace{p(t^i=1|x^i)}_{\alpha} = p(y^i=1|x^i)$

e) Ok... In other words, you are giving me $h(x^i)$ and assuming $p(t^i=1, x^i) = \text{either } 0 \text{ or } 1$, depending on x^i (ie noiseless model), and you wish me to compute $\langle h(x^i), x^i | y^i=1 \rangle$ (probability of "correct prediction").

Well, we know for a fact that there are no "false positives" here, so $h(x) = 0$ if $t^i = 0$. $h(x) | t^i=1 = \alpha$ from 2d.

$$\begin{aligned} \Rightarrow E[h(x^i) | y^i=1] &= E[p(y^i=1|x^i) | y^i=1] = E[\underbrace{\{p(y^i=1|t^i=1, x^i)p(t^i=1, x^i)}_{=0} + \\ &\quad \underbrace{\{p(y^i=1|t^i=0, x^i)p(t^i=0, x^i)\}}_{=1} | y^i=1] = E[\alpha p(t^i=1|x^i) | y^i=1] \\ &= \alpha E[p(t^i=1 | y^i=1, x^i)] = \underline{\alpha} \end{aligned}$$

2f) Done.

3a) Exponential-family distributions have pdfs of the form $p(y; \eta) = b(y) e^{T(y) - a(\eta)}$.
The Poisson-distribution is an example of these, with:
 $b(s) = 1/y!$, $T(y) = y$, $a(\eta) = \log \lambda$, $a(\eta) = \lambda = \exp(\eta)$.

b) $J(\eta) = E[y; \eta] = \langle p(y; \lambda) \rangle = \lambda = \exp(\eta)$

~~Theta~~ ~~theta~~ ~~theta~~ ~~theta~~

$$\begin{aligned}
 &= \Theta^T x^i \\
 &= \exp(\eta) = \exp(\Theta^T x^i) = h_\Theta(x^i) \\
 3c) \quad J(\Theta) &= \sum_{i=1}^n \log p(y^i | x^i; \Theta) = \sum_{i=1}^n \log b(y^i) + \underbrace{\eta(x^i; \Theta)}_{\text{learning rate}, > 0.} y^i - \underbrace{a(\eta(x^i; \Theta))}_{\cancel{\Theta_j + \alpha \sum_{j=1}^n \dots}} \\
 &\Rightarrow \frac{\partial l}{\partial \Theta_j} = \sum_{i=1}^n x_j^i y^i - h_\Theta(x^i) x_j^i \\
 &\Rightarrow \Theta_j := \Theta_j + \alpha \sum_{i=1}^n [y^i - h_\Theta(x^i)] x_j^i \Leftrightarrow \Theta := \Theta + \alpha \sum_{i=1}^n [y^i - h_\Theta(x^i)] x^i
 \end{aligned}$$

3d) Done.

$$4a) \frac{\partial}{\partial \eta} \int p(y; \eta) dy = 0 = \int \frac{\partial p}{\partial \eta} dy = \int p(y; \eta) \cdot \frac{\partial [\eta y - a(\eta)]}{\partial \eta} dy$$

$$= \int p(y - \frac{\partial a}{\partial \eta}) dy \Leftrightarrow \underline{\langle y \rangle} = E[\underline{Y; \eta}] = \underline{\frac{\partial a}{\partial \eta}}$$

$$4b) \frac{\partial^2}{\partial \eta^2} \int p dy = 0 = \int \frac{\partial^2}{\partial \eta^2} \left[p(y - \frac{\partial a}{\partial \eta}) \right] dy = \int \frac{\partial^2 p}{\partial \eta^2} (y - \frac{\partial a}{\partial \eta})$$

$$+ p(-\frac{\partial^2 a}{\partial \eta^2}) dy \Leftrightarrow \underline{\frac{\partial^2 a}{\partial \eta^2}} = \int p(y - \frac{\partial a}{\partial \eta})^2 dy = \underline{\langle y^2 \rangle} - \underline{\langle y \rangle}^2 = \underline{\text{Var}[Y; \eta]}$$

$$4c) l(\Theta) = -\sum_{i=1}^n \log p(y^i | x^i; \Theta) = \sum_{i=1}^n \log b(y^i) + \underbrace{\eta(x^i; \Theta)}_{\Theta^T x^i} y^i - \underbrace{a(\eta(x^i; \Theta))}_{\Theta^T x^i}$$

$$\Rightarrow \frac{\partial l}{\partial \Theta_j} = -\sum_{i=1}^n x_j^i y^i - \frac{\partial a}{\partial \eta} x_j^i = -\sum_{i=1}^n \left(y^i - \frac{\partial a}{\partial \eta} \right) x_j^i \Rightarrow \frac{\partial^2 l}{\partial \Theta_k \partial \Theta_j} = \sum_{i=1}^n \frac{\partial^2 a}{\partial \eta^2} x_k^i x_j^i \Leftrightarrow \nabla_\Theta^2 l = \sum_{i=1}^n \frac{\partial^2 a}{\partial \eta^2} x^i x^i$$

$$\text{Now, } \underline{w^T \nabla_\Theta^2 l w} = \sum_{i=1}^n \frac{\geq 0}{\frac{\partial^2 a}{\partial \eta^2}(x^i; \Theta)} (w^T x^i)^2$$

$$= \sum_{i=1}^n \underbrace{\text{Var}(Y|x^i; \Theta)}_{\geq 0 \text{ always}} \cdot \underbrace{\langle w, x^i \rangle^2}_{\geq 0 \text{ always}} \geq 0 \quad \forall w \Rightarrow \nabla_\Theta^2 l \text{ is psd and thus } l(\Theta) \text{ is}$$

a convex function, making it easy to minimise.

5a)

$$\text{log-lik} \downarrow$$

$$1) \underline{\ell(\theta)} = \sum_{i=1}^n \log p(y^i | \hat{x}^i; \theta) = \sum_{i=1}^n \underbrace{+ \log(2\pi)^{1/2} \sigma^{-1}}_{= \theta^T \hat{x}^i} + \frac{1}{2\sigma^2} [y^i - h_\theta(\hat{x}^i)]^2$$

$$2) \frac{\partial \ell}{\partial \theta_j} = \sum_{i=1}^n \frac{1}{2\sigma^2} (-2\hat{x}_j^i y^i + 2h_\theta(\hat{x}^i) \hat{x}_j^i) = \underbrace{-\frac{1}{\sigma^2} \sum_{i=1}^n \hat{x}_j^i (y^i - h_\theta(\hat{x}^i))}_{= \theta^T \hat{x}^i}$$

$$\Rightarrow \underline{\theta_j := \theta_j + \alpha \sum_{i=1}^n (y^i - h_\theta(\hat{x}^i)) \hat{x}_j^i} \Leftrightarrow \underline{\theta = \theta + \alpha \sum_{i=1}^n (y^i - h_\theta(\hat{x}^i)) \hat{x}^i}$$

5b) Done.

5c) ~~Using~~ Taylor expansion, it is possible to train fit a sinusoid (which is what the training data is, along with some added noise) when the number of polynomials is sufficient.

With this in mind, we see from the plot that $k=1, 2, 3$ is too few polynomials, 5 & 10 is sufficient, while 20 is too many (~~the model is fitting the noise~~ the model is ~~fitting~~ fitting the model noise).

5d) Now (obviously) the $k=1$ model fits the best, as one of its features because it contains the correct feature $\sin(x)$, which matches the structure of the data perfectly, and little other features, making overfitting a non-issue. Predictably, adding more higher-order polynomials increases overfitting.

~~This fails for k=10~~ This result is in strong contrast to 5c, where ~~this~~ a longish Taylor series is needed to correctly capture the fine-wave (albeit not too long, as otherwise overfitting occurs).

PS: It's very important that the added sine feature has the same wavenumber as the one in the ~~s~~ which generated the data.

We here see extreme overfitting: $k \geq 5$ fits exactly to the datapoints.

5e)