

1) Logistic Regression: Training Stability.

(a) What is the most notable difference in training the logistic regression model on datasets A & B?

\Rightarrow Training on dataset A converged to epsilon almost immediately, while dataset B does not seem to converge ever.

b) Investigate why the training procedure behaves unexpectedly on dataset B, but not on A.

\Rightarrow When graphing the two datasets; we see that dataset B is linearly separable; (i.e. there \exists a hyperplane such that separates the two classes with zero error).

- This is an issue for the MLE for logistic regression because:

$$\hat{\theta}_{MLE} = \underset{\theta}{\operatorname{argmax}} \ell(\theta)$$

$$\text{where } \ell(\theta) = - \sum_{i=1}^m \log \left(1 + e^{-y^i \cdot \theta^T x^i} \right)$$

- You want to $\operatorname{argmax} \ell(\theta)$; since $\ell(\theta)$ is negated; you want the maximum, meaning as least negative or possible (smallest negative value)

\Rightarrow so you want

$$\log(1 + e^{-y^i \cdot \theta^T x^i}) \rightarrow 0$$

to be as small as possible; so approach 0.

- looking specifically, you want $\log(\text{value})$ to be small as possible; so you want value to be as small as possible.

so you want $1 + e^{-y^i \cdot \theta^T x^i} \rightarrow 0$

so you want $e^{-y^i \cdot \theta^T x^i} \rightarrow 0 \dots$

- looking at this; since its $e^{-\text{exponent}}$; you want exponent $\rightarrow \infty$ for $e^{-y^i \cdot \theta^T x^i}$ to approach 0

• so objective is essentially $y^i \cdot \theta^T x^i \rightarrow \infty$

\Rightarrow DATASET B case (linearly separable)

- If a hyperplane when your classification error = 0.
- All training examples $\Rightarrow \{x^i, y^i | y^i \cdot \theta^T x^i > 0\}$ (because all correctly classified)

• so given $y^i \cdot \theta^T x^i > 0 \forall (x^i, y^i)$ to make $y^i \cdot \theta^T x^i$ we can blow up θ : make it bigger & bigger \rightarrow so as $\theta \rightarrow \infty$, our MLE continues to drop.

- hence; we never converge; as simply $\uparrow \theta$; our performance on MLE; gets better.

Versus Dataset A (not linearly separable)

- \nexists a hyperplane where classification error = 0
 - \exists points where $\{x^i, y^i | y \Rightarrow y^i \cdot \theta^T x^i < 0\}$
 - So there are points (x^i, y^i) such that $\sum y^i \cdot \theta^T x^i$ is $\sum -y^i \cdot \theta^T x^i$ i.e., decrease

$$Y = \theta^T X$$

- What does this mean? As $\theta \uparrow$; these points start "decrease" $Y \cdot \theta^T X$; get bigger as well; as θ is bigger
 - $\therefore \uparrow \theta$ eventually "hurts" performance on MLE as \exists a θ ; such that $Y \cdot \theta^T X$ is maximized where $\theta \neq \infty$.
 - So we have a point of convergence; no more improvement fr..

- TLDR's Dataset A is not linearly separable; so it converges because there is a max value of MLE; where $\theta \neq \infty$.

Dataset B is linearly separable; so the Max value of MLE is $\theta = \infty$; so θ keeps getting pushed larger \Rightarrow to infinity & does not converge.

1c) For each of these possible modifications, state whether or not it would lead to the algo converging on datasets, like B.

i. Using a different constant - learning rate

\Rightarrow No, the key issue is that the stopping condition is

$$\|\theta_{\text{old}} - \theta_{\text{new}}\| < e^{-15}$$

✓

• depends on $\theta_{\text{new}} = \theta_{\text{old}} - \text{learning rate} * \text{grad.}$

• grad won't change; but theoretically learning rate could be small enough such that our $\Delta\theta$ hits the stopping condition. But the issue is that for this to happen the condition would be hit immediately in the first iteration since θ is constant, so there would be no change to converge.

ii. Decreasing the learning rate over time (i.e. scaling the initial learning rate by $1/t$, where t is the number of grad descent iterations thus far).

Yes this would work, we would reach an optimal solution; but it would eventually stop because learning rate would be small enough to trigger stopping condition

✓

iii. Linear Scaling of the input features



⇒ No. The resulting points will still be linearly separable.

iv. Adding a regularization term $\|\theta\|_2^2$ to the loss function



Yes. This would penalize making θ infinitely big, eventually forcing a convergence

v. Adding a zero-mean Gaussian noise to the training data or labels

Yes, adding noise could make the data not linearly separable.



Q) Are support vector machines, which use the hinge loss, vulnerable to datasets like B? Why or why not? Give an informal justification

Quick intermission on hinge loss.

Hinge loss: For each sample (x^i, y^i) , where $y^i \in \{-1, 1\}$: Hinge Loss = $\max(0, 1 - y^i \cdot w^T x^i)$

• You either get 0; or $(1 - y^i \cdot w^T x^i)$ as your value.

• This is loss, you want to keep it as low as possible

Correct prediction:

$w^T x^i$ = same sign as y^i

$$\therefore y^i \cdot w^T x^i > 0 \quad (\text{IS margin } \geq 1)$$

$$\therefore (1 - y^i \cdot w^T x^i) \leq 0 \quad (\text{as well})$$

$$\therefore \text{Hinge loss} = \max(0, 1 - y^i \cdot w^T x^i) = 0$$

so for correct predicting you get hinge loss = 0

Incorrect prediction:

$w^T x^i \neq$ same sign as y^i

$$\therefore y^i \cdot w^T x^i < 0 \quad (\text{or is margin } \leq 1)$$

$$\therefore (1 - y^i \cdot w^T x^i) > 0 \quad (\text{as well})$$

$$\therefore \text{Hinge loss} = \max(0, 1 - y^i \cdot w^T x^i) = (\text{+}) \text{ positive number}$$

so for incorrect predicting you get hinge loss = pos number = you get punished.

SVM minimizes

$$\min \underbrace{\frac{1}{2} \|w\|^2}_{\text{keeps the margin}} + C \cdot \underbrace{\sum_{i=1}^m \text{Hinge Loss}}_{\text{minimize classification error}}$$

keeps the margin
large
(regularization)

minimize classification error
(using hinge loss)

No. This is because hinge loss penalizes wrong predictions (adding a positive loss), but keeps hinge loss as 0 for correct classification. There is no incentive to push W larger, as the loss doesn't get better.

\hookrightarrow the problem is the "Theta to infinity problem"

\Rightarrow hinge loss = $\max(0, 1 - y \cdot \hat{y})$, where $\hat{y} = w^T x + b$

• so with lively sympathy

$$y - \hat{y} > 0 \quad (\text{since prediction is always right})$$

so hinge loss = $\max(0, \text{negative number})$

so hinge loss = 0; and objective

function $J(\hat{y}) = 0$ (no incentive to make w bigger).

2) Model calibration:

- Will try to understand the output $h_\theta(x)$ of the hypothesis function of a logistic regression model, in particular what we might treat the output as a probability.
- When the probabilities output by a model match empirical observation, the model is "well-calibrated".
- ex; a well calibrated model's for $S \in \{x^i\}$ for which $h_\theta(x^i) \approx 0.7$, around 70% of those example should have positive labels. In a well-calibrated model; this is true at every probability value
- Log reg tends to output well-calibrated probabilities; often not true for other classifiers like naive Bayes or SVM's
- Suppose we have a training set $\{(x^i, y^i)\}_{i=1}^n$ with $x^i \in \mathbb{R}^{n+1}$ and $y^i \in \{0, 1\}$. Assume we have an intercept term $x_0^i = 1$ for all i . Let $\theta \in \mathbb{R}^{n+1}$ be the maximum likelihood parameters learned after fitting a log regression model. In order for the model to be considered well-calibrated, given any range probability (a, b) such that $0 \leq a \leq b \leq 1$ and any examples x^i where the model outputs $h_\theta(x^i)$ fall in the range (a, b) the fraction of positives in that set of examples should be equal to the average of the model outputs for those examples. That is $\frac{1}{n} \sum h_\theta(x^i)$ must hold

$$\frac{\sum_{c \in I_{a,b}} p(y^c = 1 | x^c; \theta)}{|c \in I_{a,b}|} = \frac{\sum_{c \in I_{a,b}} \mathbb{I}_{\{y^c = 1\}}}{|\{c \in I_{a,b}\}|}$$

$$P(y=1|x; \theta) = h_\theta(x) = \frac{1}{1 + \exp(-\theta^\top x)}$$

$I_{a,b} = \{c | c \in \{1, \dots, m\}, h_\theta(x^c) \in (a, b)\}$ is

an index set of training examples x^c where $h_\theta(x^c) \in (a, b)$ and $|S|$ denotes the size of these S .

a) Show that the above property holds true for the described log reg model over the range $(a, b) = (0, 1)$

Hint use the fact that we include a bias term.

\Rightarrow MLE for logistic regression:

\Rightarrow Log reg \circ

$$f(\theta) = \frac{1}{1 + e^{-\theta^\top x}} = \frac{1}{1 + \exp(-\theta^\top x)}$$

MLE \circ

$$s(y=1|x; \theta) = h_\theta(x)$$

$$s(y=0|x; \theta) = 1 - h_\theta(x)$$

$$\Rightarrow \text{since } L(\theta) = \prod_{i=1}^m f(x^i | \theta)$$

$$\Rightarrow L(\theta) = \prod_{i=1}^m h_\theta(x^i) \cdot \prod_{i=1}^m (1 - h_\theta(x^i))$$

$$\Rightarrow L(\theta) = \prod_{i=1}^m h_\theta(x^i)^{y^i} \cdot \prod_{i=1}^m (1 - h_\theta(x^i))^{1-y^i}$$

$$\Rightarrow l(\theta) = \sum_{i=1}^m y^i \log h(x^i) + (1-y^i) \log (1 - h(x^i))$$

• MLE occurs when $\nabla l(\theta) = 0$

$$0 = \nabla l(\theta) = (y - h_\theta(x)) x_j$$

from
← prev
derivation
(notes!)

$$0 = (y - h_\theta(x)) x_j$$

$\Rightarrow x_j = j^{\text{th}} \text{ feature. Let } j = 0 \text{ (0}^{\text{th}} \text{ feature).}$

• given bias term $x_0^i = 1$; plug in

$$\Rightarrow 0 = (y - h_\theta(x)) \mid : y^i = \mathbb{I}\{y^i = 1\}$$

$$0 = y - h_\theta(x)$$

$$y = h_\theta(x)$$

$$\therefore h_\theta(x^i) = P(y^i = 1 | x^i; \theta)$$

\therefore equivalent statement
(plug in)

$$\sum_{i \in I_{0,1}} \mathbb{I}\{y^i = 1\} = \sum_{i \in I_{0,1}} P(y^i = 1 | x^i; \theta)$$

- Divide both by the same val $\{ i \in I_{a,b} \}$

$$= \frac{\sum_{i \in I_{0,1}} \mathbb{I}\{y(i) = 1\}}{\sum_{i \in I_{0,1}} 1} \stackrel{def}{=} \frac{\sum_{i \in I_{0,1}} p(y(i) = 1 | x(i); \theta)}{\sum_{i \in I_{0,1}} 1}$$

2b) If we have a binary classification model that is perfectly calibrated — that is, the property we just proved holds for any $(a, b) \in [0, 1]^2$ — does this necessarily imply that the model achieves perfect accuracy? Is the converse necessarily true? Justify your answers.

\Rightarrow No, since it only guarantees that the fraction of all total answers are the same prediction as l_0 belief; it can still be the case that individual values of x_i don't match.

* Crash course on calibration

Well-calibrated confidence matches reality.

Poorly calibrated: Confidence is misguided.

Calibration means that the model's confidence score reflects the true likelihood of correctness.

- * A model can be overfit or under-fit
- * A well calibrated model \Rightarrow Accuracy \approx confidence

* Hence the converse \Rightarrow A perfectly accurate model does not mean its well calibrated; because the model can still have confidence that doesn't match its accuracy, I.e., 100% accuracy; but the model is only 80% confident. A well calibrated model with 80% confidence would have $\geq 80\%$ accuracy.

2c) Discuss what effect including L_2 regularization in the logistic regression objective has on model calibration.

\Rightarrow The L_2 regularization you add to break such

$$l(\theta) = \sum_{i=1}^m y^i \log h(x^i) + (1-y^i) \log(1-h(x^i)) + \lambda \|\theta\|_2^2$$

$$\nabla l(\theta) = \sum_{i=1}^m (y - h_\theta(x^i)) x^i + \nabla \lambda \|\theta\|_2^2$$

$$\gamma_j = 1$$

$$\nabla l(\theta) = 0 = \sum_{i=1}^m y - h_\theta(x^i) + \nabla \lambda \left(\sum_{j=1}^n \theta_j^2 \right)$$

$$\Rightarrow \lambda \frac{\partial}{\partial \theta} \sum_{j=1}^n \theta_j^2$$

$$\Rightarrow \lambda \sum_{j=1}^n 2\theta_j$$

$$0 = \sum_{i=1}^m y_i - h_\theta(x^{(i)}) + \gamma \sum_{j=1}^n 2\theta_j$$

$$\sum_{i=1}^m h_\theta(x^{(i)}) - \gamma \sum_{j=1}^n 2\theta_j = \sum_{i=1}^m y_i$$

- no longer fits well-calibrated property

$$\frac{\sum_{i \in I_{a,b}} P(y_i = 1 | x^{(i)}; \theta)}{|\{i \in I_{a,b}\}|} = \frac{\sum_{i \in I_{a,b}} \prod_{j \neq i} (y_j - 1)}{|\{i \in I_{a,b}\}|}$$

∴ It makes the logistic regression model no longer well-calibrated.

* Note: When logistic regression (GLM assumption) holds, all ranges $(a, b) \subset [0, 1]$ are well-calibrated.

In addition when S_{train} & S_{test} are from the same distribution and when the model has not overfit or underfit, logistic regression tends to be well-calibrated on unseen test data as well. It is therefore good in practice; when we care about level of uncertainty in model output.

3) Bayesian interpretation of regularization

- MAP estimation is diff from estimation technique we learned so far; like MLE because it allows us to incorporate prior knowledge into the estimate,

a) Show that $\theta_{\text{MAP}} = \arg \max_{\theta} P(y|x; \theta) p(\theta)$ if we assume that $p(\theta) = p(\theta|x)$

$$\theta_{\text{MAP}} = \arg \max_{\theta} P(\theta|x, y)$$

$$\Rightarrow P(\theta|x, y) = \frac{P(y|x, \theta) P(\theta|x)}{P(y|x)}$$

$$\theta_{\text{MAP}} = \arg \max_{\theta} = \frac{P(y|x, \theta) P(\theta|x)}{P(y|x)}$$

$\arg \max_{\theta}$ \Rightarrow find the val of θ that maximizes

$$\therefore \arg \max_{\theta} \frac{P(y|x, \theta) P(\theta|x)}{P(y|x)}$$

$$= \underset{\theta}{\operatorname{argmax}} \ p(y|x; \theta) p(\theta|x)$$

(since $p(y|x)$ is independent of θ)

$$\Rightarrow p(\theta) = p(\theta|x)$$

↓ sub in

$$\therefore \underset{\theta}{\operatorname{argmax}} \ p(y|x; \theta) p(\theta|x)$$

$$\Rightarrow \therefore \underset{\theta}{\operatorname{argmax}} \ p(y|x; \theta) p(\theta)$$

b) Recall that L_2 regularization penalizes the L_2 norm of the parameters while minimizing the loss (i.e. negative log likelihood in case of probabilistic models). Now we will show that MAP estimation with a zero-mean Gaussian prior over θ , specifically $\theta \sim \mathcal{N}(0, \gamma^2 I)$ is equivalent to applying L_2 regularization with MLE estimation. Spec. show that.

$$\theta_{\text{MAP}} = \underset{\theta}{\operatorname{argmin}} -\log p(y|x, \theta) + \gamma \|\theta\|_2^2$$

$$\theta_{\text{MAP}} = \underset{\theta}{\operatorname{argmax}} p(y|x; \theta) p(\theta)$$

$$= \underset{\theta}{\operatorname{argmin}} -p(y|x; \theta) p(\theta)$$

$$= \underset{\theta}{\operatorname{argmin}} -[\log p(y|x; \theta) + \log p(\theta)]$$

$$\theta_{\text{MAP}} = \underset{\theta}{\operatorname{argmin}} -\log p(y|x, \theta) - \log p(\theta)$$

= $\gamma \|\theta\|_2^2$ in MLE

$$\text{try: } \log P(\theta) \text{ given } \theta \sim N(0, \eta^2 I)$$

PDF of a d -dimensional multivariate normal distribution $N(\mu, \Sigma)$ for a vector x is given by:

$$P(x) = \frac{1}{(2\pi)^{\frac{d}{2}} |\Sigma|^{\frac{1}{2}}} \exp\left(-\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu)\right)$$

$$P(\theta) = \frac{1}{(2\pi)^{\frac{d}{2}} |\Sigma|^{\frac{1}{2}}} \exp\left(-\frac{1}{2} (\theta - \mu)^T \Sigma^{-1} (\theta - \mu)\right)$$

$$= P(\theta) = \frac{1}{(2\pi)^{\frac{n}{2}} |\Sigma|^{\frac{1}{2}}} \exp\left(-\frac{1}{2} (\theta)^T \Sigma^{-1} (\theta)\right)$$

$$= \frac{1}{2\pi^{\frac{n}{2}} |\eta^2 I|^{\frac{1}{2}}} \exp\left(-\frac{1}{2} (\theta)^2 (\eta^2 I)^{-1}\right)$$

$$= (2\pi)^{-\frac{n}{2}} (\eta^2 I)^{-\frac{1}{2}} \exp\left(-\frac{1}{2} (\theta)^2 (\eta^2 I)^{-1}\right)$$

\Rightarrow sub into θ_{MAP}

$$\theta_{MAP} = \arg \min_{\theta} \left\{ -\log P(y|x; \theta) - \log P(\theta) \right\}$$

$$\begin{aligned} \Theta_{MAP} &= \underset{\theta}{\operatorname{argmin}} \left[-\log p(y|x;\theta) - \frac{1}{2} \log \left[(2\pi)^{\frac{n}{2}} (\eta^2 I)^{\frac{1}{2}} \exp\left(-\frac{1}{2}\theta^T(\eta^2 I)^{-1}\theta\right) \right] \right] \\ &\dots \underset{\theta}{\operatorname{argmin}} \left[\log \left[(2\pi)^{\frac{n}{2}} (\eta^2 I)^{\frac{1}{2}} \exp\left(-\frac{1}{2}\theta^T(\eta^2 I)^{-1}\theta\right) \right] \right] \\ &= \underset{\theta}{\operatorname{argmin}} \left[\log (2\pi)^{\frac{n}{2}} + \log (\eta^2 I)^{-\frac{1}{2}} + \left(-\frac{1}{2}\theta^T(\eta^2 I)^{-1}\theta\right) \right] \\ &\quad \curvearrowleft \quad \curvearrowleft \\ &\text{constants wrt } \theta. \underset{\theta}{\operatorname{argmin}} \theta = \underline{\text{drop}} \\ &= \underset{\theta}{\operatorname{argmin}} -\left[\frac{1}{2} \theta^T (\eta^2 I)^{-1} \theta \right] \\ &\quad \vdots \text{ plug back in} \\ \Theta_{MAP} &= \underset{\theta}{\operatorname{argmin}} \left(-\log p(y|x;\theta) - \left[-\frac{1}{2} \theta^T (\eta^2 I)^{-1} \theta \right] \right) \\ &= \underset{\theta}{\operatorname{argmin}} -\log p(y|x;\theta) + \frac{1}{2} \theta^T (\eta^2 I)^{-1} \theta \\ &\geq \underset{\theta}{\operatorname{argmin}} -\log p(y|x;\theta) + \frac{1}{2} \frac{\theta^2}{\eta^2} \\ \Rightarrow \text{Original Form:} &\quad \checkmark \\ \Theta_{MAP} &= \underset{\theta}{\operatorname{argmin}} -\log p(y|x;\theta) + \lambda \|\theta\|_2^2 \\ \frac{1}{2} \frac{\theta^2}{\eta^2} &= \lambda \|\theta\|_2^2 \end{aligned}$$

$$\frac{1}{2n^2} \theta^2 = \gamma \| \theta \|_2^2$$

$$\therefore \gamma = \frac{1}{2n^2}$$

c) Now consider a specific instance, a linear regression model given by $y = \theta^T x + \epsilon$; where $\epsilon \sim N(0, \sigma^2)$.

Like before, assume a Gaussian prior on this model such that $\theta \sim N(0, \eta^2 I)$. For notam, let X be the design matrix of all the training example inputs, and each row vector \vec{x} one example input, and \vec{y} be the column vector of all the example outputs. Some up with closed form for MAP

$$\theta_{MAP} = \arg \max_{\theta} p(\theta | X, y)$$

* MAP \Rightarrow probability of θ given X , y .
Let's find the most likely θ given the \vec{y} . which means both X & y !

$$P(A|B) = \frac{P(B|A) \cdot P(A)}{P(B)}$$

$$\Rightarrow \theta_{MAP} = \arg \max_{\theta} p(\theta | X, y) = \frac{p(y | X, \theta) \cdot p(\theta | X)}{p(y | X)}$$

$\arg \max_{\theta} \rightarrow p(y|x)$ is irrelevant for maximizing since
it doesn't depend on θ

$$\Rightarrow \theta_{MAP} = \arg \max_{\theta} P(y|x, \theta) \cdot P(\theta|x)$$

* Bayesian stats assumes θ is not something discrete,
it's something we believe. Even before seeing any data,
we have a prior belief about θ . θ exists
independently outside of data we observe. Hence,
"just knowing x doesn't change what we believe about
 θ ".
Hence $P(\theta|x) = P(\theta)$

$$\Rightarrow \theta_{MAP} = \arg \max_{\theta} p(y|x, \theta) \cdot P(\theta)$$

PDF for theta:

$$P(x) = \frac{1}{(2\pi)^{\frac{n}{2}} |\Sigma|^{\frac{1}{2}}} \exp \left(-\frac{1}{2} (x-\mu)^T \Sigma^{-1} (x-\mu) \right)$$

$$P(\theta) = \frac{1}{(2\pi)^{\frac{n}{2}} |\Sigma|^{\frac{1}{2}}} \exp \left(-\frac{1}{2} (\theta)^T \Sigma^{-1} (\theta) \right)$$

$p(y|x)$... multivariate form as x is also a
vector ... needed μ & Σ

$$\epsilon \sim N(0, \sigma^2)$$

$$\theta \sim N(0, \eta^2 I)$$

$$p(y|x, \theta) = ?$$

• Find μ of distribution

$$y = \theta^T x + \epsilon$$



$$\therefore y = x\theta + \epsilon = \mu$$

, covariance, $x\theta$ is deterministic, so covariance comes from ϵ
 $= \sigma^2$

$$\therefore \Sigma = \sigma^2 I$$

$$p(y|x, \theta) \sim N(x\theta, \sigma^2 I)$$

⇒ find p.d.f. of single variable normal distribution for y :

given $y \sim N(\mu, \sigma^2)$

$$p(y) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{-\frac{(y-\mu)^2}{2\sigma^2}\right\} \quad ; \text{ where } y \sim N(\mu, \sigma^2)$$

$$\Rightarrow p(y|x, \theta) \Rightarrow N(x\theta, \sigma^2 I)$$

(for each individual x^i)

$$p(y^i|x^i, \theta) \stackrel{N}{\Rightarrow} (x^i\theta, \sigma^2)$$

$$\mu = \theta^T x^i$$

$$\text{variance} = \sigma^2$$

$y = ? \Rightarrow$ univariate PDF, we calculate P

the pdf at each point $y^i \dots$ so $y = y^i$

$$p(y^i) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{-\frac{(y^i - \theta^T x^i)^2}{2\sigma^2}\right\}$$

$p(y^i)$ is a single point

$$p(y|x, \theta) = \prod_{i=1}^{n=1} \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{-\frac{(y^i - \theta^T x^i)^2}{2\sigma^2}\right\}$$

$$p(y|x, \theta) = \frac{1}{(\sqrt{2\pi\sigma^2})^m} \exp\left\{-\frac{1}{2\sigma^2} \sum_{i=1}^m (y^i - \theta^T x^i)^2\right\}$$

$$= \frac{1}{2\pi\sigma^2} \exp\left\{-\frac{1}{2\sigma^2} (y^i - \theta^T x^i)^2\right\}$$

$$\Rightarrow \|x\|_2 = \sqrt{x_1^2 + x_2^2}$$

$$\|\vec{y} - X\theta\|_2 = \sqrt{(y^i - \theta^T x^i)^2}$$

$$\underline{\Rightarrow \|\vec{y} - X\theta\|_2^2 = (y^i - \theta^T x^i)^2}$$

$$p(y|x, \theta) = \frac{1}{2\pi\sigma^2} \exp\left\{-\frac{1}{2\sigma^2} \|\vec{y} - X\theta\|_2^2\right\}$$

$$\text{From before... } \log p(\theta) = \frac{1}{2\eta^2} \|\theta\|_2^2$$

\Rightarrow plug back in...

$$\theta_{MAP} = \arg \min_{\theta} -\log p(y|x, \theta) - \log p(\theta)$$

$$\theta_{MAP} = \arg \min_{\theta} -\log p(y|x, \theta) + \frac{1}{2\eta^2} \|\theta\|_2^2$$

$$\theta_{MAP} = \arg \min_{\theta} \left[-\log \left[\frac{1}{2\pi^{\frac{m}{2}} \sigma^m} \exp \left(-\frac{1}{2\sigma^2} \|\vec{y} - X\theta\|_2^2 \right) \right] \right] \dots$$

$$= \arg \min_{\theta} \left(-\log \frac{1}{2\pi^{\frac{m}{2}} \sigma^m} - \frac{1}{2\sigma^2} \|\vec{y} - X\theta\|_2^2 \right) \dots$$

drop b/c $\arg \min_{\theta}$ (constnt)

$$= \arg \min_{\theta} \left(-\frac{1}{2\sigma^2} \|\vec{y} - X\theta\|_2^2 - \frac{1}{2\eta^2} \|\theta\|_2^2 \right)$$

$$\therefore J(\theta) = -\frac{1}{2\sigma^2} \|\vec{y} - X\theta\|_2^2 - \frac{1}{2\eta^2} \|\theta\|_2^2$$

$$= \frac{1}{2\sigma^2} \|X\theta - \vec{y}\|_2^2 - \frac{1}{2\eta^2} \|\theta\|_2^2$$

$$\Rightarrow \|X\| = \sqrt{x^T x}$$

$$\|x\|^2 = x^T x \quad (\text{u-sub})$$

$$= \frac{1}{2\sigma^2} \nabla(x^T x - \vec{y}^T x) - \frac{1}{2\eta^2} \|\theta\|_2^2$$

$$= \frac{1}{2\sigma^2} \nabla (x\theta - \bar{y})^T (x\theta - \bar{y}) - \frac{1}{2\eta^2} \|\theta\|_2^2$$

gradient of quadratic sum

$$= 2x^T(x\theta - \bar{y})$$

$$= \frac{1}{\sigma^2} \cdot 2x^T(x\theta - \bar{y}) + \frac{1}{2\eta^2} \|\theta\|_2^2$$

$$= \frac{1}{\sigma^2} x^T(x\theta - \bar{y}) + \frac{1}{\eta^2} \cdot \theta^T \theta$$

$$\nabla J = \frac{1}{\sigma^2} (x^T x\theta - x^T \bar{y}) + \frac{1}{\eta^2} \theta$$

$$\nabla J_\theta = 0$$

$$\Rightarrow 0 = \frac{1}{\sigma^2} (x^T x\theta - x^T \bar{y}) + \frac{1}{\eta^2} \theta$$

$$0 = \frac{1}{\sigma^2} x^T x\theta - \frac{1}{\sigma^2} x^T \bar{y} + \frac{1}{\eta^2} \theta$$

$$\frac{1}{\sigma^2} x^T \bar{y} = \frac{1}{\sigma^2} x^T x\theta + \frac{1}{\eta^2} \theta$$

$$\frac{1}{\sigma^2} x^T \bar{y} = \theta \left(\frac{1}{\sigma^2} x^T x + \frac{1}{\eta^2} \right)$$

$$\frac{\frac{1}{\sigma^2} x^T \bar{y}}{\frac{1}{\sigma^2} x^T x + \frac{1}{\eta^2}} = \theta$$

$$x^T \tilde{y} \cdot \left(x^T x + \frac{\sigma^2}{\eta^2} \right)^{-1} = \theta$$

$$\theta_{MAP} = \arg \min_{\theta} J(\theta) = x^T \tilde{y} \cdot \left(x^T x + \frac{\sigma^2}{\eta^2} \right)^{-1}$$

c) Laplace

$$f_L(z|\mu, b) = \frac{1}{2b} \exp\left(-\frac{|z-\mu|}{b}\right)$$

- consider a linear regression model given by $y = x^T \theta + \varepsilon$
- $\varepsilon \sim N(0, \sigma^2)$. Assume a laplace prior $\theta \sim L(0, bI)$

• Find θ_{MAP}

$$\theta_{MAP} = \arg \min \theta - \log p(y|x, \theta) - \log p(\theta)$$

$$f_L(\theta) = \frac{1}{2b} \exp\left(-\frac{|\theta - \mu|}{b}\right)$$

$$\theta = (z|\mu, b)$$

$$\theta = z, \mu = 0$$

$$f_L(\theta) = \frac{1}{2b} \exp\left(-\frac{|\theta^i - 0|}{b}\right)$$

$$f_L(\theta) = \prod_{i=1}^{i=n} \frac{1}{2b} \exp\left(-\frac{|\theta^i|}{b}\right)$$

Since ~~and~~ symmetry

$$f_L(y|x, \theta) = ?$$

$$\Rightarrow y = x^\top \theta + \epsilon ; \text{ where } \epsilon \sim N(0, \sigma^2)$$

$$y \sim N(0, \sigma^2)$$

$$f_L(y|x, \theta) = \prod \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(y^i - \theta^\top x^i)^2}{2\sigma^2}\right)$$

\Rightarrow

$$\theta_{MAP} = \arg \min_{\theta} -\log p(y|x, \theta) - \log p(\theta)$$

$$= \arg \min_{\theta} -\log \left[\frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(y^i - \theta^\top x^i)^2}{2\sigma^2}\right) \right] - \log \left[\left(\frac{1}{2\sigma} \right)^n \exp\left(-\frac{\|\theta\|^2}{2\sigma^2}\right) \right]$$

$$= \arg \min_{\theta} \sum_{i=1}^m \frac{(y^i - \theta^\top x^i)^2}{2\sigma^2} + \sum_{i=1}^n \frac{1}{b} |\theta_i|$$

$$= \arg \min_{\theta} \frac{1}{2\sigma^2} \sum_{i=1}^m (y^i - \theta^\top x^i)^2 + \sum_{i=1}^n \frac{1}{b} |\theta_i|$$

$$\|x\|_2 = \sqrt{\sum x^i} = \sqrt{\sum x^i}^2 = \|x\|_2^2$$

$$= \arg \min_{\theta} \frac{1}{2\sigma^2} \|x\theta - \bar{y}\|_2^2 + \frac{1}{b} \|\theta\|_1$$

$$= \arg \min_{\theta} \|x\theta - \bar{y}\|_2^2 + \frac{2\sigma^2}{b} \|\theta\|_1$$

$$\boxed{\therefore J(\theta) = \|x\theta - \bar{y}\|_2^2 + \frac{2\sigma^2}{b} \|\theta\|_1}$$

$$\boxed{\theta_{MAP} = \arg \min_{\theta} J(\theta)}$$

• what is γ^2

$$\Rightarrow J(\theta) = \|\mathbf{x}\theta - \vec{y}\|_2^2 + \gamma \|\theta\|,$$

$$J(\theta) = \|\mathbf{x}\theta - \vec{y}\|_2^2 + \frac{\gamma^2}{b} \|\theta\|^2,$$

$\therefore \gamma = \frac{\gamma^2}{b}$

4) Constructing Kernels:

⇒ First: kernels walkthrough

1) What is a kernel?

⇒ A kernel such $K(x, z)$ is a way to compute a dot product between feature maps without explicitly mapping inputs to a high dimensional space. That is, you may map or map

$$\phi : \mathbb{R}^n \rightarrow \mathcal{H}$$

\mathcal{H} is some Hilbert space (higher dimensions); and the kernel function is defined using this mapping ϕ

$$K(x, z) = \langle \phi(x), \phi(z) \rangle$$

• This allows kernel-based algos; like SVMs to work in \mathcal{H} without ever computing $\phi(x)$ directly.

Next we need to know a valid kernel.

• A function $K(x, z)$ is a valid kernel iff :

i) Symmetric: $K(x, z) = K(z, x)$

2) positive semi definite: PSD : for any finite set $\{x^1, \dots, x^m\}$
 the matrix $K \in \mathbb{R}^{m \times m}$ with entries
 $K_{ij} = k(x^i, x^j)$ satisfies
 $\forall c \in \mathbb{R}^m$, $(^T K c \geq 0)$ (all eigenvalues of kernel matrix K are non-negative)

* Core Idea: A kernel is a symmetric, positive semi-definite function. Any operation that preserves symmetry + positive semi-definiteness gives another valid kernel

Kernel construction rules

• If k_1 and k_2 are valid kernels then:

1) Addition: $k(x, z) = k_1(x, z) + k_2(x, z) \Rightarrow$ valid kernel

2) Scalar multiplication: $k(x, z) = a k_1(x, z)$ for $a > 0$ is valid

3) Multiplication: $k(x, z) = k_1(x, z) k_2(x, z) \Rightarrow$ valid

4) Function application: If f is a function and $k(x, z) = f(x)f(z)$ then this is a valid kernel because $\phi(x) = f(x)$, a 1-D feature map

5) Composition: If $\phi: \mathbb{R}^n \rightarrow \mathbb{R}^d$, and k_3 is a kernel over $\mathbb{R}^d \times \mathbb{R}^d$, then $k(x, z) = k_3(\phi(x), \phi(z))$ is also a kernel

6) polynomials with positive coefficients: If $p(x) \mapsto$ positive coefficients, and k_1 is a kernel

then $\Phi(K, (x, z))$ is also a kernel.

Finding counterexamples \Rightarrow prove that it's not PSD?

• Using the quadratic form \Rightarrow a such $x^T A x$ that takes a vector x and outputs a scalar value, where A is a symmetric matrix and x^T represents the transpose x .

• K is not positive semi-definite if you can find a vector z , such that

$$\underline{z^T K z < 0}$$

• for example

$$K = K_1 - K_2$$

$$\text{let } K_1 = 2K_2$$

$$\text{then } K_2 = \frac{1}{2}K_1$$

$$\text{then: } K = K_1 - 2K_2 = -K_1$$

and therefore
 $\underline{z^T K z = z^T -K_1 z \leq 0} \dots$

• therefore $z^T K z < 0$; and K is not positive semi-definite

* Remember that positive semi-definite dfn is for a matrix K , its quadratic form $z^T K z \geq 0$.

Back to Q4.

k_1 & k_2 are kernels $\mathbb{R}^n \times \mathbb{R}^n$, & $a \in \mathbb{R}^+$
& $f: \mathbb{R}^n \mapsto$ be a real valued funcn; let $\phi: \mathbb{R}^n \mapsto \mathbb{R}^d$
be a funcn mapping from \mathbb{R}^n to \mathbb{R}^d . Let k_3 be
a kernel over $\mathbb{R}^d \times \mathbb{R}^d$, and let $p(x)$ be a
polynomial over X with positive coefficients -

(a) $k(x, z) = k_1(x, z) + k_2(x, z)$

k_1 is PSD & k_2 is PSD

$k_1 + k_2$ is also PSD

$$k = k_1 + k_2$$

PSD defn:

$$z^T k z \geq 0$$

$$\Rightarrow z^T (k_1 + k_2) z$$

$$(z^T k_1 + z^T k_2) z$$

$$\underline{z^T k_1 z + z^T k_2 z \geq 0}; \text{ so } k = k_1 + k_2 \text{ is PSD}$$

+ a valid kernel

b) $k(x, z) = k_1(x, z) - k_2(x, z)$

$$\Rightarrow z^T k z \geq 0$$

$$\text{let } k_2 = 2k_1$$

$$\Rightarrow k = k_1 - k_2 = k_1 - 2k_1 = -k_1$$

$$\Rightarrow z^T k z = z^T -k_1 z = -z^T k_1 z$$

$$\therefore -z^T k_1 z < 0$$

$$\therefore z^T k z < 0$$

$\therefore k = k_1 - k_2$ is not PSD & not a valid kernel

c) $k(x, z) = a k_1(x, z)$; $a \in \mathbb{R}^+$ (positive real num)

PSD defn: $z^T k z \geq 0$

$$z^T k z = z^T a k_1 z$$

$$= a z^T k_1 z$$

k_1 is a valid kernel; $\therefore z^T k_1 z \geq 0$
 a is a positive real number: $a > 0$

$$\therefore a z^T k_1 z = \text{pos}$$

$\oplus \cdot \oplus$

$$\therefore a z^T k_1 z \geq 0$$

$\therefore k = a k_1$ is PSD & a valid kernel

d) $k(x, z) = -a k_1(x, z)$

PSD defn: $z^T k z \geq 0$

$$K = -\alpha k_1(x, z)$$

$$z^T K z = z^T -\alpha k_1 z = -\alpha z^T k_1 z$$

k_1 is PSD; so $z^T k_1 z \geq 0$

$$\therefore -\alpha z^T k_1 z < 0$$

$\ominus \times \oplus$

$$\therefore z^T K z < 0$$

$\therefore K = -\alpha k_1$ is not PSD & not a valid kernel

e) $K(x, z) = k_1(x, z) k_2(x, z)$

PSD defn: $z^T K z \geq 0$

$$z^T K z = z^T k_1 k_2 z$$

• Mercer's theorem, every valid kernel is an inner product in some (possibly infinite-dimensional) space

$$K(x, x') = \langle \phi(x), \phi(x') \rangle$$

x , x' are in your original input space

$\phi(x), \phi(x')$ are the higher dimensional feature values, given by a ϕ (feature map)

• $K(x, x') =$ the inner product of $\phi(x)$ and $\phi(x')$

- $K(x, x')$ (Kernel function) gives you the inner product between $\phi(x)$ & $\phi(x')$ without ever the need to compute $\phi(x)$

Mercer's theorem

If $K(x, x')$ is a symmetric, positive definite kernel, then it can be written as: $K(x, x') = \sum_{i=1}^{\infty} \lambda_i \phi_i(x) \phi_i(x')$

Where:

- $\lambda_i \geq 0$ are eigenvalues
- $\phi_i(x)$ are orthogonal eigenfunctions
- This sum acts like a dot product in a (possibly infinite-dimensional) space
- So it gives us a feature mapping

$$\phi(x) = (\sqrt{\lambda_1} \phi_1(x), \sqrt{\lambda_2} \phi_2(x), \dots)$$

$$\text{then } K(x, x') = \langle \phi(x), \phi(x') \rangle$$

- So Mercer guarantees that a valid kernel function K , is an inner product in some feature space.
- guarantees the existence of a mapping $\phi(x)$.
- Kernel trick works b/c "you can trust this kernel" — it's

the dot product of some high-dimensional represen~~ts~~, even if you never~~ts~~ it

$$k_1(x, z) = \langle \phi(x), \phi(z) \rangle$$

$$k_2(x, z) = \langle \phi_2(x), \phi_2(z) \rangle$$

(From Mercer's theorem)

ϕ & ϕ_2 are w^g functions that create a M dimensional & N dimensional vector (higher featurespace)

• write the product $k_1 + k_2$ in terms of $\phi_1 + \phi_2$.

$$k = k_1 \cdot k_2$$

$$k(x, z) = \langle \phi(x), \phi(z) \rangle \cdot \langle \phi_2(x), \phi_2(z) \rangle$$

$$= \sum_{i=1}^m \phi(x^i) \phi(z^i) \cdot \sum_{i=1}^n \phi_2(x^i) \phi_2(z^i)$$

$$= \sum_{i=1}^m \sum_{j=1}^n [\phi(x^i) \phi_2(x^i)] [\phi(z^i) \phi_2(z^i)]$$

$$\text{let } c(x^i) = \phi(x^i) \phi_2(x^i)$$

$$\hat{=} \sum_{i=1}^m \sum_{j=1}^n c(x^i) (z^i)$$

$$= \underline{C(X)^T C(Z)}$$

C is a (m, n) dimensional vector; s.t. $C = \phi_1(x) \phi_2(x)$

Since ϕ is a feature map C ; that

can map values x to $C(x)$; and we
can denote $K(x, z)$ as a inner product;
 $C(X)^T C(Z)$; by Mercer's theorem, $K(x, z)$ is
a valid kernel.

$$f) K(x, z) = f(x) f(z)$$

$\Rightarrow f: \mathbb{R}^n \rightarrow \mathbb{R}$ be a real-valued function

$$\text{PSD} \Leftrightarrow z^T K z \geq 0$$

$$z^T f(x) f(z) \geq 0$$

$f(x) \Rightarrow \mathbb{R}^n \rightarrow \mathbb{R}$; so inner product

$$K(x, z) = \sum_{i=1}^m \sum_{j=1}^n f(x)_i f(z)_j$$

$$= \underline{C(X)^T C(Z)}$$

• we can write $K(x, z)$ as an inner product w.r.t. feature
map C , we know K_P is a valid kernel.

$$g) K(x, z) = k_3(\phi(x), \phi(z))$$

→ given: k_3 is a valid kernel, K is a valid kernel as $K = k_3$, and $k_3(x, z)$ is a kernel no matter the inputs.

$$h) K(x, z) = p(k_1(x, z))$$

$p(x)$ is a polynomial over x , with pos coefficients.

$$K(x, z) = p(k_1(x, z))$$

$$\Rightarrow z^T K(x, z) z \geq 0$$

$$z^T K(x, z) z = z^T p(k_1(x, z)) z$$

• $k_1(x, z)$ is a valid kernel, $k_1(x, z) \geq 0$
+ p is pos only coefficient ..

$$z^T [\oplus \cdot \oplus] z \quad (\text{pos coeff times } \geq 0 \text{ kernel func})$$

$$\therefore z^T p(k_1(x, z)) z \geq 0$$

∴ $K = p(k_1(x, z))$ is a valid kernel

$$\text{as } z^T p(k_1(x, z)) z \geq 0$$

5) Kernelizing the Perceptron.

Rmbr: perceptron :

inputs $X = [x_1, x_2 \dots x_n]$

Applies weights $\theta = [\theta_1, \theta_2 \dots]$
bias = b

- Computes weighted sum :

$$z = \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 \dots$$

- Applies activation function

$$\text{output} = \begin{cases} 1 & \text{if } z > 0 \\ 0 & \text{otherwise} \end{cases}$$

Kernel trick:

- we replace X with $\phi(x)$ (we don't know $\phi(x)$)
- But even w/out we have inner product

$$\phi(x_i) \cdot \phi(x_j) \Rightarrow k(x_i, x_j)$$

- For each misclassified point, we update the prediction

Kernelizing the perceptron

1) how to represent the weights?

- we can't store w directly; since we don't know ϕ

$$w = \sum_i a_i y_i \phi(x_i)$$

2) How to represent the new predict function

- plug back into decision function:

$$f(x) = \text{sign}(w \cdot \phi(x))$$

$\phi(x)$ normally its $w \cdot x$; but its
because we are mapping x to higher dimensional...

- put w back into it

$$f(x) = \text{sign}\left(\sum_i a_i y_i \phi(x_i) \cdot \underline{\phi(x)}\right)$$

3) we don't know ϕ , what do we do?

apply the kernel trick

$$k(x_i, x_j) = \phi(x_i) \cdot \phi(x_j)$$

- sub in the kernel wherever possible

$$f(x) = \text{Sign} \left(\sum_i a_i y_i \phi(x_i) \cdot \phi(x) \right)$$

$K(x_i, x)$

$$f(x) = \text{sign} \left(\sum_i a_i y_i \cdot K(x_i, x) \right)$$

a is a list of points that is a record of mistakes.

* we need to product $s(x_i)$ from update with a .

* To update "a" to include mistakes, set $a_i = 0$ for correct; and $a_i = 1$ for incorrect.

- How to add a learning rate?

* Make $\alpha = \text{learning rate}$ on updates for a_i . Since we want to multiply the effect of every training point by α learning rate, we can do this by setting $a_i = \alpha$ since a_i is multiplied for every point in the predict function's

$$f(x) = \text{sign} \left(\sum_i a_i y_i \phi(x_i) \cdot \phi(x) \right)$$

a)
 i. How will you implicitly represent the high-dimensional parameter vector Θ^* , including how the initial value $\Theta^0 = 0$ is represented?

By keeping track & adding all previous misclassified outputs $y_i \phi(x_i)$ in the predict function, we get a predict function

$$f(x) = \sum_{i=1}^m a_i y_i \phi(x_i) \phi(x).$$

By the kernel rule, this simplifies to

$f(x) = \sum_{i=1}^m a_i y_i K(x_i, x)$; allowing us to represent Θ^* , a vector whose dimension is the same as the feature vector $\phi(x)$ implicitly without knowing ϕ ; by using the kernel function in the predict function.

ii) How will you efficiently make a prediction on a new input x^{i+1} i.e. how will you compute $h_\Theta(x^{i+1}) = g(\Theta^{i^T} \phi(x^{i+1}))$, using your representation of Θ^* ?

\Rightarrow we take our prediction function that implicitly represents ϕ , by including the kernel function:

$$f(x) = \sum_{i=1}^m a_i y_i k(x_i, x)$$

and plug in x^{i+1} ; and apply $g(x)$ to it: so:
 $g(f(x^{i+1}))$

iii) Update rule update

- For misclassified values we add to the predict function, by add the term:

$$\frac{a_i y_i \phi(x_i)}{2}$$

a_i will be the learning rate; so that for misclassified points we add $y_i \phi(x_i)$ scaled by $\frac{1}{2}$ to "move" our future predictions.

This results in:

$$f(x) = \sum_{i=1}^m a_i y_i k(x_i, x)$$

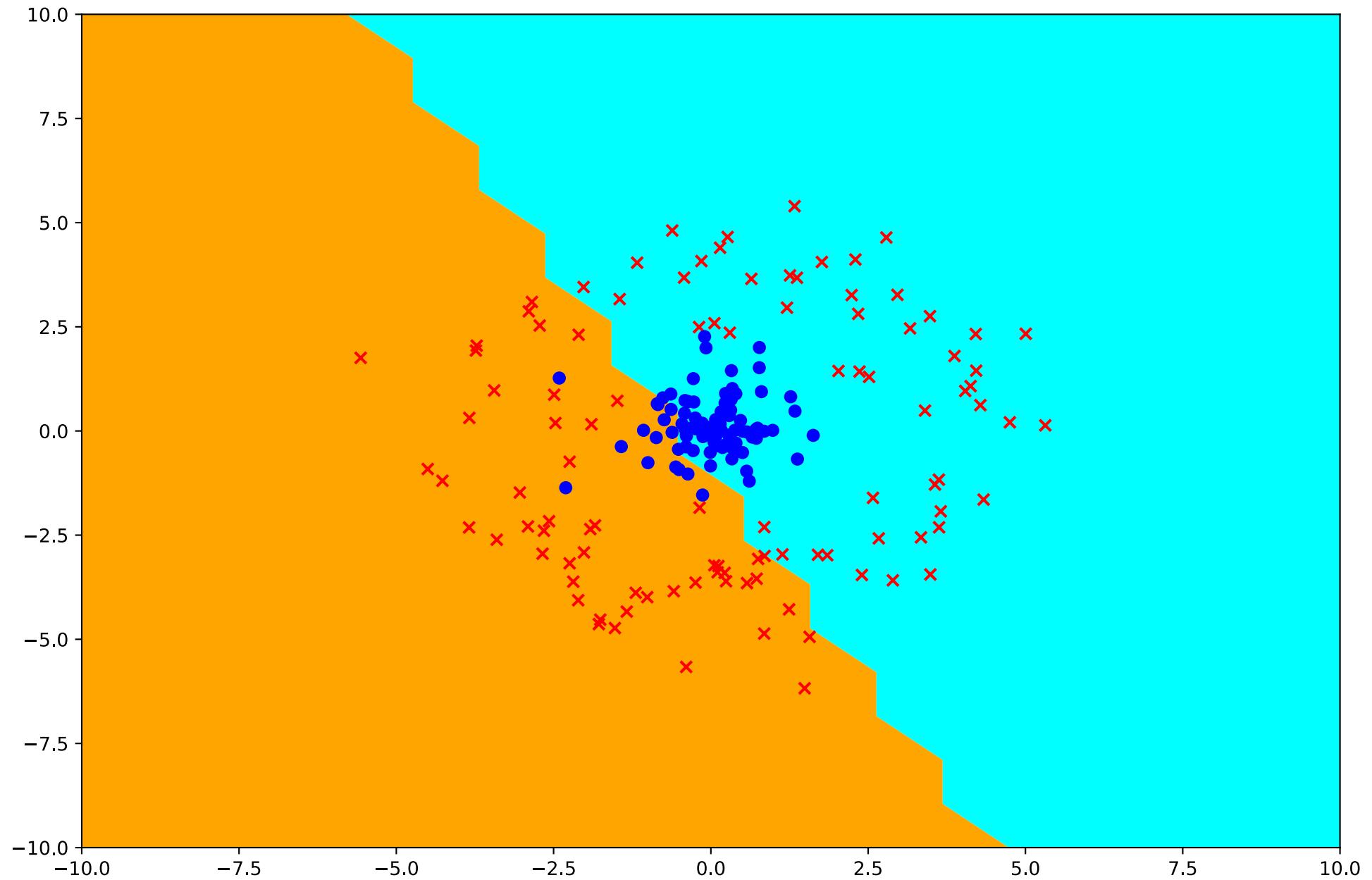
b) In PS2/src in jupyter notebook

c)

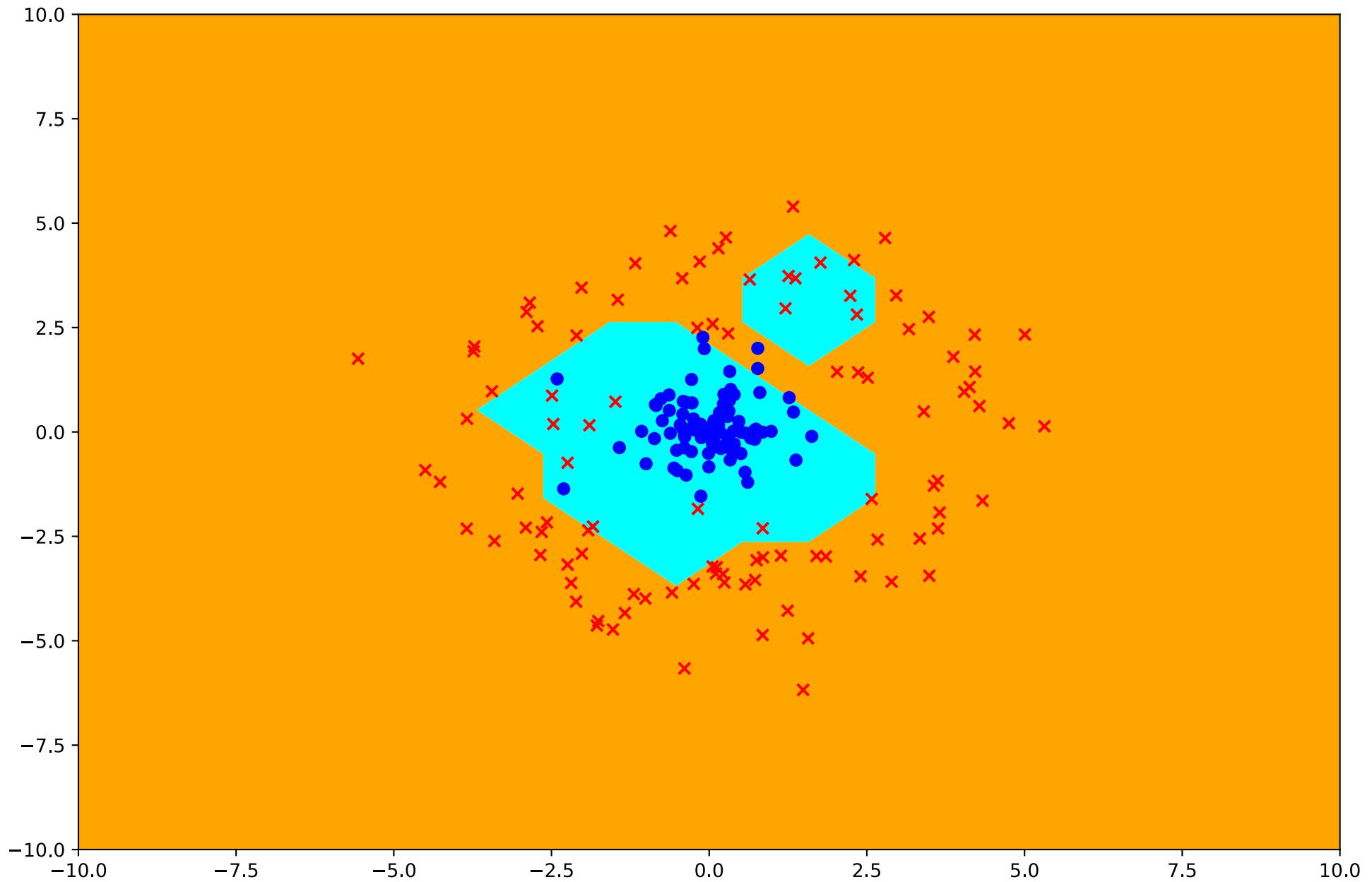
Outputs



Dot product kernel



RBF Kernel



The dot-product kernel performs badly. This is because the data is not linearly-separable, and the dot-product (linear) kernel struggles with complex, non-linear data.