

1Q.

The exponential family is a class of distributions with the following form: $p(y; \eta) = b(y) \exp(\eta^T T(y) - a(\eta))$

By identifying the parameters of the Poisson distribution:

$$p(y; \lambda) = \frac{e^{-\lambda} \lambda^y}{y!}$$

$$\exp \left[\log \left[e^{-\lambda} \lambda^y \frac{1}{y!} \right] \right]$$

$$\exp \left[-\lambda + y \log(\lambda) + \log \frac{1}{y!} \right]$$

$$\frac{1}{y!} \exp [y \log(\lambda) - \lambda]$$

so we can see: $\eta = \log(\lambda)$; $\lambda = e^\eta$; $b(y) = \frac{1}{y!}$; $T(y) = y$;
 $a(\eta) = e^\eta$.

1.6

The canonical response function for the Poisson distribution is given by:

$$g(\eta) = E(y; \eta) = \lambda = e^\eta$$

This result can also be obtained by the use of the exponential family!

$$g(\eta) = \frac{\partial}{\partial \eta} a(\eta) = \frac{\partial}{\partial \eta} e^\eta = e^\eta$$

1.C

The log-likelihood of a training example $(x^{(i)}, y^{(i)})$ is given by $\log p(y^{(i)}|x^{(i)}; \theta)$.

To derive the stochastic gradient descent update rule, we can start by computing the partial derivative of the log-likelihood with respect to parameter θ_j , with $g = \theta^T x^{(i)}$:

$$\begin{aligned}\frac{\partial}{\partial \theta_j} \log p(y^{(i)}|x^{(i)}; \theta) &= \frac{\partial}{\partial \theta_j} ((\theta^T x^{(i)})^T y^{(i)} - e^{\theta^T x^{(i)}}) \\ &= (y^{(i)} - e^{\theta^T x^{(i)}}) x_j^{(i)}\end{aligned}$$

The yields the update rule for stochastic gradient descent with learning rate λ :

$$\theta_j := \theta_j + \lambda (y^{(i)} - e^{\theta^T x^{(i)}}) x_j^{(i)}$$

2.a.

$$K(x, z) = K_1(x, z) + K(x, z)$$

is a kernel because the sum of two PSD matrices is PSD.

2.b.

$$K(x, z) = K_1(x, z) - K(x, z)$$

is not necessarily a kernel. Counterexample:

$K_1(x, z) = 0$ and $K_2(x, z) = 1 \forall x, z$ are kernels, but $K(x, z) = -1$ is obviously not.

2.c

$$K(x, z) = \alpha K_1(x, z)$$

is a kernel, because the set of PSD matrices is closed under multiplication with positive real numbers

L.d.

$$K(x, z) = -\alpha K_1(x, z)$$

would only be a kernel if:
 $\alpha \leq 0$ or $K_1 = 0$.

3. ai

The algorithm only makes one pass through the training set, it will always be the case that $\theta^{(i)}$ is a linear combination of $\phi(x^{(1)}), \dots, \phi(x^{(i)})$ i.e.

$$\theta^{(i)} = \sum_{k=1}^i \theta_k \phi(x^{(k)})$$

for some real numbers $\theta_k \in \mathbb{R}$.

Thus we could represent $\theta^{(i)}$ by the list $[\theta_1, \dots, \theta_i]$. For technical reason we will represent $\theta^{(i)}$ as pair $([\theta_1, \dots, \theta_i], [x^{(1)}, \dots, x^{(i)}])$.

In particular $\theta=0$ will be represented by $([], [])$, i.e pair

of two empty lists.

3. All

For any vector $x \in \mathbb{R}^n$ and any vector θ with

$$\theta = \sum_{k=1}^i \theta_k \phi(x^{(k)})$$

We've got:

$$\begin{aligned}\theta^\top \phi(x) &= \sum_{k=1}^i \theta_k \phi(x^{(k)})^\top \phi(x) \\ &= \sum_{k=1}^i \theta_k K(x^{(k)}, x).\end{aligned}$$

For a new input x^{i+1} we can complete the prediction

efficiently by

$$\begin{aligned} \ell_{Q^{(i)}}(x^{(i+1)}) &= g(\theta^{(i)}{}^T \phi(x^{(i+1)})) \\ &= g\left(\sum_{k=1}^i Q_k^{(i)} K(x^{(k)}, x^{(i+1)})\right). \end{aligned}$$

3. a iii'

The new update rule

$$\theta^{(i+1)} := \theta^{(i)} + \lambda (y^{(i+1)} - h_{\theta^{(i)}}(x^{(i+1)})) \phi(x^{(i+1)})$$

can be derived by appending

$$\theta_{i+1} := \lambda (y^{(i+1)} - h_{\theta^{(i)}}(x^{(i+1)}))$$

resp. $x^{(i+1)}$ to the first resp. second component of our representation $([\theta_1, \dots, \theta_i], [x^{(1)}, \dots, x^{(i)}])$ of $\theta^{(i)}$.

B.C.

With the dot product kernel, the perceptron performs poorly because the model attempts to identify a linear decision boundary but the dataset is far from linearly separable.