

$$1.1) P(y|x, w) = \mathcal{N}(\vec{w}^T \vec{x}, \sigma) = \frac{1}{\sqrt{2\pi\sigma}} e^{\left[-\frac{(y - \vec{w}^T \vec{x})^2}{2\sigma}\right]}$$

$$1.2) P(w|D) = \frac{P(w) \cdot P(D|w)}{P(D)} = \frac{P(w)}{P(D)} \cdot \prod_{i=0}^N \mathcal{N}(\vec{w}^T \vec{x}_i, \sigma)$$

$$1.3) \hat{w} = \arg \max_w P(\vec{w}|D)$$

$$= \arg \min_w -\log P(\vec{w}|D)$$

$$= \arg \min_w -\log \frac{P(\vec{w}) \cdot P(D|\vec{w})}{P(D)}$$

$$= \arg \min_w (-\log P(\vec{w}) - \log P(D|\vec{w}) + \cancel{\log P(D)})$$

\* log of the products  
we calculated in 1.2  
becomes a summation

in 1.3

$$= \arg \min_w \left( -\log \prod_{i=0}^N \mathcal{N}(\vec{w}^T \vec{x}_i, \sigma) - \log(P(\vec{w})) \right)$$

$$= \arg \min_w \left( -\sum_{i=0}^N \left( \log(\cancel{\frac{1}{\sqrt{2\pi\sigma}}}) + \log e^{-\frac{(y - \vec{w}^T \vec{x})^2}{2\sigma}} \right) - \log(P(\vec{w})) \right)$$

$$= \arg \min_w \left( -\sum_{i=0}^N \frac{(y - \vec{w}^T \vec{x})^2}{2\sigma} - \log(P(w)) \right)$$

$$= \arg \min_w \left( \cancel{\frac{1}{2\sigma}} \sum_{i=0}^N (y - \vec{w}^T \vec{x})^2 - \log \mathcal{N}(0, \alpha \mathbf{I}) \right)$$

$$= \arg \min_w \left( \sum_{i=0}^N (y - \vec{w}^T \vec{x})^2 - \log \left( \cancel{\frac{1}{\sqrt{2\pi\alpha\mathbf{I}}}} \cdot e^{\frac{1}{2\alpha\mathbf{I}} \vec{w}^T \vec{w}} \right) \right)$$

$$= \arg \min_w \left( \sum_{i=0}^N (y - \vec{w}^T \vec{x})^2 - \frac{1}{2\alpha\mathbf{I}} \vec{w}^T \vec{w} \right)$$

$$\approx \min \left( \sum (\vec{w}^T \vec{x} - y)^2 - \lambda \|\vec{w}\|_2^2 \right) \text{ which is RLR}$$

1.4) because  $p(D)$  drops out of min/max equations,

$$\arg \max_w p(D|w) = \arg \min_w -\log p(D|w) = \arg \min_w -\log p(w|D) - p(w)$$
$$= \arg \min_w \left( \sum (\eta - \vec{w}^T \vec{x})^2 \right)$$

which is the linear regression equation without the  $p(w)$  which is the regularization term so therefore this is the unregularized linear regression equation.

$$2.1) \quad \vec{w}_{k+1}^T \vec{w}_{opt} \geq \vec{w}_k^T \vec{w}_{opt} + \gamma \|\vec{w}_{opt}\|$$

$$(\vec{w}_{k+1} - \vec{w}_k)^T \vec{w}_{opt} \geq \gamma \|\vec{w}_{opt}\|$$

$$(\vec{w}_k + y_k \vec{x}_k - \vec{w}_k)^T \vec{w}_{opt} \geq \gamma \|\vec{w}_{opt}\|$$

$$(y_k x_k)^T \vec{w}_{opt} \geq \gamma \|\vec{w}_{opt}\|$$

Since  $\gamma$  is the smallest distance of any data point from the separating plane  $w$ , we know that:

$$|(y_k x_k)^T \vec{w}_{opt}| \geq \gamma \|\vec{w}_{opt}\|$$

and we know the left hand side of our equation will always be positive because  $y_k$  and  $x_k$  will always have the same sign and  $w_{opt}$  will always be positive, so we can then prove

because the classes are linearly separable

$$(y_k x_k)^T \vec{w}_{opt} \geq \gamma \|\vec{w}_{opt}\|$$

2.2

$$\begin{aligned} \|\vec{w}_{k+1}\|^2 &= \vec{w}_{k+1}^T \vec{w}_{k+1} = (\vec{w}_k + \gamma_k x_k)^T (\vec{w}_k + \gamma_k x_k) \\ &= \vec{w}_k^T \vec{w}_k + 2 \gamma_k w_k x_k + \gamma_k^2 x_k^2 \end{aligned}$$

$$\cancel{\|\vec{w}_k\|^2} + 2 \gamma_k w_k x_k + \gamma_k^2 x_k^2 \leq \cancel{\|\vec{w}_k\|^2} + 1$$

~~$$\|\vec{w}_k\|^2 + 2 \gamma_k w_k x_k + \gamma_k^2 x_k^2 \leq \|\vec{w}_k\|^2 + 1$$~~

$$2 \gamma_k w_k x_k + \gamma_k^2 x_k^2 \leq 1$$

↑  
This has to  
be negative based  
on the problem  
definition

↑  
This is  
in range  
{-1,1} so at  
worst it will  
be  $|x|^2 = 1$

↑  
 $x^2$  is normalized  
to equal 1

in the end we have a negative number  
plus the product of 2 things that each  
max at 1 so their product also maxes at  
which means the whole left side maxes  
out at less than or equal to 1

$$2.3) \quad \|w_{k+1}\| \leq \sqrt{M} \quad \text{RHS}$$

$$\|w_{k+1}\|^2 \leq M$$

$$\|w_k\|^2 + 1 \leq M$$

$$\|w_0\|^2 = 0$$

$$\|w_1\|^2 = 1$$

$$\|w_2\|^2 = 2$$

$$\|w_3\|^2 = 3$$

from the result of 2.2,  $\|w_k\|^2$  will increment by 1 each time if we start with weights of zero then  $\|w_k\|^2 + 1$  has to be less than or equal to the number of mistakes.

$$w_{opt}^T M \cdot \min_x \frac{w_{opt}^T x}{w_{opt}^T} \leq \|w_{k+1}\| \cdot w_{opt}$$

$$M \cdot \min_x w_{opt}^T x \leq \|w_{k+1}\| w_{opt} \leq$$

$$M \cdot \min_x w_{opt}^T x \leq w_k^T w_{opt} + \gamma \|w_{opt}\|$$

$$M \cdot \min_x w_{opt}^T x \leq w_k^T w_{opt} + \min \frac{w_{opt}^T x}{w_{opt}^T} \cdot w_{opt}$$

2.3 (cont')

$$w_{k+1}^T w_{opt} \geq w_k^T w_{opt} + \gamma \|w_{opt}\|$$

$$\|w_{k+1}\| \cancel{\|w_{opt}\|} \geq \|w_k\| \cancel{\|w_{opt}\|} + \gamma \cancel{\|w_{opt}\|}$$

$$\|w_{k+1}\| - \|w_k\| \geq \gamma$$

telescoping sum:  $w_{k+1} - w_k + w_k - w_{k-1} + w_{k-1} - w_{k-2} \dots w_1 - w_0$

$$\|w_{k+1}\| - \underbrace{\|w_0\|}_0 \geq \gamma$$

$$\|w_{k+1}\| \geq \gamma \quad \text{LHS}$$

$$2.4) \quad M \times \gamma \leq \|w_{k+1}\| \leq \sqrt{M}$$

$$M \times \gamma \leq \sqrt{M}$$

$$M^2 \times \gamma^2 \leq M$$

$$M^2 \leq \frac{M}{\gamma^2}$$

$$\frac{M^2}{M} \leq \gamma^{-2}$$

$$M \leq \gamma^{-2}$$

2.5) if 2 classes are not linearly separable, this would cause 2.1's conclusion to not be provable which would mean we also could not prove 2.3 so we could not prove that the perceptron converges

$$3) \quad \hat{w} = \arg \max_w \sum_{i=1}^N y_i \vec{w}^T \vec{x}_i + \lambda (\vec{w}^T \vec{w} - 1)$$

$$\hat{w} = \arg \max_w \sum_{i=1}^N y_i \vec{w}^T \vec{x}_i + \lambda (\|w\|^2 - 1)$$

$$\frac{df}{dw} = 0$$

$$\frac{df}{dw} = \sum_{i=1}^N y_i x_i + 2\lambda |\vec{w}| = 0$$

$$\sum_{i=1}^N y_i x_i = -2\lambda |\vec{w}|$$

$$\frac{\sum_{i=1}^N y_i x_i}{-2\lambda} = \hat{w}$$

$$\sum_{i=1}^N y_i x_i \propto \hat{w}$$

$$\hat{w} \propto \sum_{i: x_i \in \mathcal{C}_1} y_i x_i + \sum_{j: x_j \in \mathcal{C}_{-1}} y_j x_j$$

$$\hat{w} \propto \sum_{i: x_i \in \mathcal{C}_1} x_i - \sum_{j: x_j \in \mathcal{C}_{-1}} x_j$$