

Risk: % you will get it wrong. Expected value of loss.

$$\log p_{\theta}(y|x) = N(f_{\theta}(x), \Sigma_{\theta}(x)) = -\frac{1}{2} (f_{\theta}(x) - y)^T \Sigma_{\theta}(x)^{-1} (f_{\theta}(x) - y) - \frac{1}{2} \log |\Sigma_{\theta}(x)| + \text{const}$$

$$(L(\theta, x, y)) \quad \text{if } \Sigma_{\theta}(x) = I \rightarrow -\frac{1}{2} \|f_{\theta}(x) - y\|^2 + \text{const} \approx \text{MSE}$$

overfitting - low emp. risk, high risk

underfitting - high emp. risk, high risk

expected value of error w.r.t. data dist (theorectically): $E_{\theta \sim p(\theta)} [\|f_{\theta}(x) - f(x)\|^2]$

"how many will our predictions be, on average."

$$\text{let } \tilde{f}(x) = E_{\theta \sim p(\theta)} [f_{\theta}(x)], \quad E_{\theta \sim p(\theta)} [\|f_{\theta}(x) - f(x)\|^2]$$

$$= E_{\theta \sim p(\theta)} [\|f_{\theta}(x) - \tilde{f}(x) + (\tilde{f}(x) - f(x))\|^2]$$

$$= \underbrace{E_{\theta \sim p(\theta)} [\|f_{\theta}(x) - \tilde{f}(x)\|^2]}_{\text{Variance}} + \underbrace{E_{\theta \sim p(\theta)} [\|\tilde{f}(x) - f(x)\|^2]}_{\|\tilde{f}(x) - f(x)\|^2 \rightarrow \text{Bias}^2}$$

"how much does our prediction change w/ this dataset" error never goes away. "Bias + Variance Tradeoff"

Regularization: something we add to the loss function to reduce variance.

Bayesian perspective: Given \mathcal{D} , what is most likely θ ?

$$p(\theta|\mathcal{D}) = \frac{p(\mathcal{D}, \theta)}{p(\mathcal{D})} \propto p(\mathcal{D}, \theta) = p(\mathcal{D}|\theta) p(\theta) \leftarrow \text{prior, how likely } \theta \text{ is before } \mathcal{D}.$$

$$\text{New loss: } -\left(\sum_{i=1}^N \log p(y_i|x_i; \theta)\right) - \log p(\theta) \leftarrow \text{choose this}$$

based on $N(0, \sigma^2)$.

Regularized Lin Reg: $\log p(\theta) \propto -\lambda \|\theta\|^2 + \text{const}$ ($\lambda = \frac{1}{2\sigma^2}$), same for Reg. Log. Reg.

L1 Reg. MAE - sparsity

L2 Reg. MSE - smooth

Dropout: used for NNS

Numerical perspective: regularizer makes undetermined problems well-determined.

Optimization Perspective: regularizer makes loss landscape easier to search.

Regularizer introduces hyperparameters that we have to select for success.

ML Workflow

Train	- optimize to find θ , hyperparams.
Valid.	- tune hyperparameters, select model class, which features to use
Test	- "final exam", report final performance.