

Artificial Neural Networks (Gerstner). Exercises for week 4

Regularization and Tricks of the Trade

Exercise 1. Softmax function for multi-class problems.

We have a network with K output units indexed by $1 \leq k \leq K$. We want that output unit k gives the posterior probability for class k

$$\hat{y}_k = P(C_k|\mathbf{x}) \quad (1)$$

Show that we arrive at the softmax function if we set the joint probabilities as $p(\mathbf{x}, C_k) = e^{a_k}$.

Hints: Proceed analogously to the 'derivation' of the sigmoidal unit in class last week. Remember that The softmax function is

$$\hat{y}_k = \frac{\exp(a_k)}{\sum_j \exp(a_j)} \quad (2)$$

Exercise 2. Bagging and dropout

Assume K variants of a model, with model k having test error $E_k = E_0 + \epsilon_k$, where ϵ_k has mean $\mathbb{E}[\epsilon_k] = 0$, auto-variance $\mathbb{E}[\epsilon_k^2] = v$ and co-variance $\mathbb{E}[\epsilon_k \epsilon_n] = c$.

- What is the expected value of the *bagging test error*, i.e. the expected test error of the model obtained by averaging over all variants?
- What is the variance of the bagging test error?
- For dropout: do you expect the co-variance c to be close to zero, between zero and v , close to v or larger than v ? Justify your answer with one sentence.

Exercise 3. Bias and variance of gradient estimators

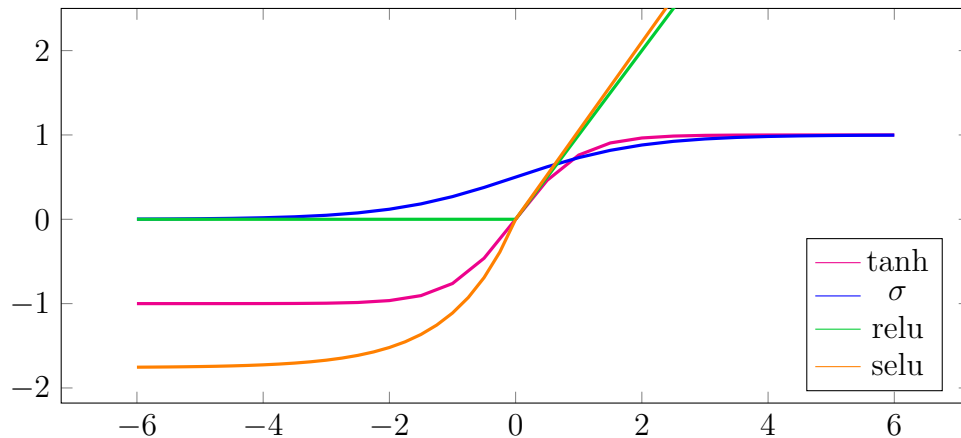
For training data $(\mathbf{x}^1, y^1), \dots, (\mathbf{x}^P, y^P)$ and some loss $E(\mathbf{w}) = \frac{1}{P} \sum_{\mu} \ell(f_{\mathbf{w}}(\mathbf{x}^{\mu}), y^{\mu})$, the gradient is given by $\nabla E(\mathbf{w}) = \frac{1}{P} \sum_{\mu} \nabla \ell(f_{\mathbf{w}}(\mathbf{x}^{\mu}), y^{\mu})$, with e.g. $\ell(x, y) = \frac{1}{2}(x - y)^2$.

- In each step of stochastic gradient descent one sample $(\mathbf{x}^{\mu}, y^{\mu})$ of the training data is selected. Show that $\nabla \ell(f_{\mathbf{w}}(\mathbf{x}^{\mu}), y^{\mu})$ is an unbiased estimator of $\nabla E(\mathbf{w})$, if each training sample is selected with equal probability. Hint: An estimator $\hat{\theta}$ of a quantity θ is called unbiased, if its expectation $\mathbb{E}[\hat{\theta}] = \theta$.
- Instead of single sample stochastic gradient descent it is common practice to use mini-batches. Show that the mini-batch gradient estimator $\frac{1}{M} \sum_{i=1}^M \nabla \ell(f_{\mathbf{w}}(\mathbf{x}^i), y^i)$, with $1 < M < P$, has lower variance than the single sample estimator, if the samples (\mathbf{x}^i, y^i) in each mini-batch are sampled uniformly from the training data.

Exercise 4. Different activation functions

The choice of the non-linearity function $g(x)$ can have a significant impact on learning speed and final performance. Which non-linearity is best, is still an active research question; the favorite non-linearity in the last century was probably the hyperbolic tangent $\tanh(x)$; since 2010, the rectified linear unit $\text{relu}(x) = \max(0, x)$ is highly popular and there is a fair chance that the new favorite will be the scaled exponential linear unit $\text{selu}(x) = \lambda x$ if $x > 0$ and

$\text{selu}(x) = \alpha(\exp(x) - 1)$ otherwise, with $\lambda \approx 1.0507$ and $\alpha \approx 1.75814$. Currently, it seems that the key concepts to discuss the different non-linearities are, first, *using the nonlinearity*, second the *vanishing gradient problem* and, third, the *bias shift problem*.



a. Using the nonlinearity

- (i) Show that a multi-layer neural network with linear activation function $g(x) = x$ is equivalent to a single layer linear network. Hint: the product of two matrices is again a matrix.
- (ii) Assume that in each layer the inputs follow a Normal distribution with mean zero and small variance, i.e. $\sigma^2 \ll 1$. For which of the activation functions $\sigma(x) = 1/(1 + \exp(-x))$, $\tanh(x)$, $\text{relu}(x)$ and $\text{selu}(x)$ is a deep network basically equivalent to a linear network for this input distribution?

b. Vanishing gradient problem

- (i) Assume now the inputs are such that they also fall into the non-linear regimes. For simplicity we assume that in each layer the inputs are $x_1 = -10, x_2 = -5, x_3 = -1, x_4 = 1, x_5 = 5, x_6 = 10$. Without a calculator, determine the fraction of values close to zero of $g(x_i)$ and $g'(x_i)$ for all i and $g = \sigma, \tanh, \text{relu}, \text{selu}$. For example, for \tanh none of the values $\tanh(-10), \tanh(-5), \dots, \tanh(10)$ is close to zero but $4/6 = 2/3$ of the values of $\tanh' = 1 - \tanh^2$ are close to zero.
- (ii) The update of a weight w_{ij} is proportional to $g'(x_i) \cdot g(x_j)$. Determine the fraction of $g'(x_i) \cdot g(x_j)$ that are close to zero considering all combinations of x_i and x_j and all activations $g = \sigma, \tanh, \text{relu}, \text{selu}$.
- (iii) The δ 's in backpropagation are in each layer multiplied with g' . Consider backpropagation through 3 layers, i.e. terms like $g'(x_i)g'(x_j)g'(x_k)$. Determine the fraction of such terms that are close to zero for $g = \sigma, \tanh, \text{relu}, \text{selu}$.

c. Bias shift problem

Consider a simple classification task. The data exist in \mathcal{R}^N . Data points from C_0 (with target $t = 0$) are uniformly distributed in each dimension such that $x_i \in [1, 2]$ for $i = 1 \dots N$. Data points from C_1 (with target $t = 1$) are uniformly distributed in each dimension such that $x_i \in [3, 4]$ for $i = 1 \dots N$. We want to learn to classify points using a logistic sigmoid unit trained with the cross-entropy loss; from last week, this results in the weight update rule

$$\Delta w_i = \eta \cdot (t - y) \cdot x_i$$

where $y = \sigma\left(\sum_i^N w_i x_i\right)$.

Points are presented one at a time (i.e. stochastic gradient descent).

- (i) Assume we start with all weights $w_i = 0$ and present the point \mathbf{x}^a from C_0 , update the weights, then present \mathbf{x}^b . Give the drive $a = \sum_i^N w_i x_i^b$ of the output unit in response to \mathbf{x}^b , in terms of η , \mathbf{x}^a and \mathbf{x}^b . Note: we do not yet need to specify which class \mathbf{x}^b belongs to.

- (ii) In general, we can encounter oscillations in stochastic gradient descent if a single training example strongly affects the network output – for instance, if it results in the same network output for any possible input.

We assume that if $a < -5$, $y \approx 0$, and if $a > 5$, $y \approx 1$. Under what conditions will the network output y be the same for all possible inputs \mathbf{x}^b after the first training step? Can we choose a small enough η to prevent this, independent of N ? What if we had chosen \mathbf{x}^a from C_1 instead?

- (iii) A common input normalization technique to to remove the mean from the dataset, such that $E[x_i] = 0$ across all dimensions x_i . Assume that each data point has an equal probability of coming from either C_0 and C_1 . What are the new data ranges for C_0 and C_1 after removing the mean? Repeating step (ii), do we get the same result?

- (iv) Consider a deep network where each hidden layer uses either the hyperbolic tangent or logistic sigmoid nonlinearity. Given what we've seen above, can you suggest one non-linearity or the other between hidden layers? Note that one layer's output is another layer's input.

- d. Summarize your results by ranking the different activation functions for each of the problems discussed in this exercise.

	linear for $x \sim \mathcal{N}(0, \sigma \ll 1)$	vanishing gradient problem	bias shift problem
tanh			
σ			
relu			
selu			