

MACHINE LEARNING: DEEP LEARNING

Eric Nalisnick
Johns Hopkins University

Last Update:
February 3, 2025

Contents

1	Supervised Learning with Univariate Linear Models	3
1.1	Predictive Modeling	3
1.2	Univariate Linear Model for Real-Valued Responses	3
1.3	Maximum Likelihood Estimation: A General Recipe	6
1.4	Generalized Linear Models	9
1.5	Univariate Logistic Regression for Binary Labels	10
1.6	Gradient Descent	11
1.7	Models of Artificial Neurons & the Perceptron	13
2	Supervised Learning with Multiple Linear Regression	15
2.1	Multiple Features and Feature Expansions	15
2.2	Revisiting Gradient Descent with Multivariate Derivatives	17
2.3	Multiple Output Dimensions	17
3	Model Evaluation, Model Selection, and Regularization	17
4	Feedforward Neural Networks	17
4.1	Adaptive Basis Functions	17
5	Convolutional Neural Networks	17
6	Recurrent Neural Networks	17
7	Transformers	17
8	Deep Generative Models	17
9	Uncertainty Estimation	17

1 Supervised Learning with Univariate Linear Models

The first topic we will discuss is predictive modeling using linear models—that is, models that are linear in their parameters. This will provide the building blocks with need to eventually stack these ‘shallow’ models into the ‘deep’ models that give this course its title.

1.1 Predictive Modeling

Consider the task of *predictive modeling*. Imagine that we are creating a system that, given a medical image, can predict if the patient has a particular disease, e.g. pneumonia. Such a system will be used by bringing patients into the clinic to perform the imaging, and then once the image is taken, the image will be passed to some sort of predictive model, that will generate the prediction that the radiologist will consider to inform their diagnosis. Let \mathcal{X} denote a feature space, which in the example above, is the space of all valid medical images. You can think of this as a matrix in which entries are the pixel values, intensities, or some other property of the image. Let \mathcal{Y} denote the label / response space, which in the above setting is a discrete encoding of the potential diseases. Our goal is to design some model $\hat{y} = f(\mathbf{x})$ that takes features $\mathbf{x} \in \mathcal{X}$ as input and outputs an accurate prediction $\hat{y} \in \mathcal{Y}$.

Data Generating Process Ideally, we want the above model $f(\mathbf{x})$ to match the true underlying process that generated the data. In the medical imaging example, this means that $f(\mathbf{x})$ would faithfully capture whatever is the underlying medical process that results in a person, with that given image, having the biological conditions that present as their true clinical diagnosis. Mathematically, we can say that the world generates these diseases according to a distribution $\mathbb{P}(y|\mathbf{x})$, and thus the goal of predictive modeling is to have $f(\mathbf{x}) = \mathbb{P}(y|\mathbf{x})$. Although, in practice, we are often satisfied with a close approximation.

Training Data As we will see below, we will construct $f(\mathbf{x})$ in a data-driven way. That is, instead of just hand-engineering rules or some other function for $f(\mathbf{x})$, we will *learn* a good predictive model from *data* that represents or encodes our problem of interest. Ideally, we would like to know and work with $\mathbb{P}(y|\mathbf{x})$ directly, but this is usually never the case in practice. And if we did have access to $\mathbb{P}(y|\mathbf{x})$, why then would we need to train a model $f(\mathbf{x})$? In practice, we usually just have samples from $\mathbb{P}(y|\mathbf{x})$, i.e. $y \sim \mathbb{P}(y|\mathbf{x})$. For the features \mathbf{x} , we will also assume we have samples from another underlying generative process $\mathbf{x} \sim \mathbb{P}(\mathbf{x})$. One could try to model $\mathbb{P}(\mathbf{x})$ in addition to $\mathbb{P}(y|\mathbf{x})$; this is usually called *generative modeling*, a topic we will get into later in the course. We will assume that, for purposes of training data, we are able to collect N samples of feature-label pairs, making our N -element training set $\mathcal{D} = \{(\mathbf{x}_n, y_n)\}_{n=1}^N$.

1.2 Univariate Linear Model for Real-Valued Responses

We will now get into our first (or many) concrete instantiations of $f(\mathbf{x})$, and we will start with a (seemingly) simple function: the line, with one slope parameter. Assume for the time being that the features $x \in \mathbb{R}$ and $y \in \mathbb{R}$ are both real-valued, unconstrained scalar variables. We will define the *univariate linear model* as $f(x; \mathbf{w}) \triangleq \mathbf{w} \cdot \mathbf{x}$, where x is a scalar feature value and \mathbf{w} is a scalar parameter that we wish to learn from data. To pick apart

the notation, $f(x; w)$ means that we have a function of the features x and the function is determined by parameters w . This model encodes a very simple predictive relationship: the prediction \hat{y} is proportional or inversely proportional to the feature value x .

Loss Function Given an N -sample training set \mathcal{D} and the linear model $f(x; w)$, the next step is to fit the model to the data. An intuitive way to do this is to define a *loss function* that quantifies how far off the model's predictions are from the observed data. While we will later give a complete recipe for deriving loss functions, one natural choice for real-valued, unconstrained data is the squared loss: $\ell(f; x, y) = (f(x) - y)^2$. Clearly, this will be zero when $f(x) = y$ and grow quadratically as $f(x)$ makes worse and worse predictions. Also notice that this loss doesn't care if the predictions under or over estimate y , which could be inappropriate for some applications. For example, in the American game show *The Price is Right*, contestants had to guess the sale price of items, and if they overestimated the price, they instantly lost. If you were building an AI agent to play *The Price is Right*, you would certainly want to train it with a loss function that treats over- and under- estimates differently. Now that we have devised a loss for one data point, we can compute the loss over the full training set by summing the losses for each data point:

$$\ell(w; \mathcal{D}) = \frac{1}{N} \sum_{n=1}^N \ell(w; x_n, y_n) = \frac{1}{N} \sum_{n=1}^N (f(x_n; w) - y_n)^2 = \frac{1}{N} \sum_{n=1}^N (w \cdot x_n - y_n)^2. \quad (1)$$

Note that these loss functions are a function of *the model*, with the data treated as a constant, because we want to assess how well the model fits the data and not vice versa.

Optimizing a Loss Function Now that we have defined a loss function, we want to use it to find the best setting of the parameter w . This boils down to the following optimization problem:

$$\begin{aligned} w^* &= \arg \min_w \ell(w; \mathcal{D}) \\ &= \arg \min_w \frac{1}{N} \sum_{n=1}^N (f(x_n; w) - y_n)^2 \\ &= \arg \min_w \frac{1}{N} \sum_{n=1}^N (w \cdot x_n - y_n)^2. \end{aligned} \quad (2)$$

Thus, w^* will be the parameter that minimizes the squared distance between the model predictions and the training responses y . It is unlikely the value of the loss will be exactly zero when computed using $f(x_n; w^*)$, so when we speak of 'minimizing the loss,' it is constrained by the best training performance achievable under the fixed model class—which in this case, is the univariate linear model. The loss function might be able to be driven to exactly zero if we were to choose a different model, especially one that can represent more flexible functions than a line.

Now how should we find the exact form of w^* . Fortunately, for linear models, we can do this exactly and in 'closed form,' meaning that we can get an explicit equation for w^* . This will not be the case for most of the course, and we'll often have to resort to approximate, numerical techniques. Yet, in all cases, we will rely upon tools from calculus.

Recall that the points at which a derivative equals zero represent the *critical points* of a function, meaning that that point can be a maxima, minima, or saddle point. For this linear model with the squared loss, fortunately there is just one (non-trivial) critical point and it represents the global minimum. While a proper course on optimization would go into the details of validating this claim, we will mostly ignore these details since deep learning methodologies often need to reply upon so many approximations that such proofs are not that informative of practice.

Moving on to the mechanics of taking the derivative of the loss with respect to the model parameter, we have:

$$\begin{aligned}
\frac{d}{dw} \ell(w; \mathcal{D}) &= \frac{d}{dw} \left[\frac{1}{N} \sum_{n=1}^N (w \cdot x_n - y_n)^2 \right] \\
&= \frac{1}{N} \sum_{n=1}^N \frac{d}{dw} [(w \cdot x_n - y_n)^2] \\
&= \frac{1}{N} \sum_{n=1}^N 2 \cdot (w \cdot x_n - y_n) \cdot \frac{d}{dw} [w \cdot x_n - y_n] \\
&= \frac{1}{N} \sum_{n=1}^N 2 \cdot (w \cdot x_n - y_n) \cdot x_n \\
&= \frac{2}{N} \left\{ \left(\sum_{n=1}^N w \cdot x_n^2 \right) - \left(\sum_{n=1}^N y_n \cdot x_n \right) \right\}.
\end{aligned} \tag{3}$$

Now we can find w^* by setting the derivative to zero and solving for w :

$$\begin{aligned}
0 &= \frac{d}{dw} \ell(w; \mathcal{D}) = \frac{2}{N} \left\{ \left(\sum_{n=1}^N w \cdot x_n^2 \right) - \left(\sum_{n=1}^N y_n \cdot x_n \right) \right\} \\
\Rightarrow 0 &= \left(\sum_{n=1}^N w \cdot x_n^2 \right) - \left(\sum_{n=1}^N y_n \cdot x_n \right) \\
\Rightarrow \sum_{n=1}^N w \cdot x_n^2 &= \sum_{n=1}^N y_n \cdot x_n \\
\Rightarrow w &= \frac{\sum_{n=1}^N y_n \cdot x_n}{\sum_{n=1}^N x_n^2} \triangleq w^*.
\end{aligned} \tag{4}$$

We have finally arrived at the ‘trained’ version of our model: computing $\sum_{n=1}^N y_n \cdot x_n / \sum_{n=1}^N x_n^2$ will give the value that we should plug in for the optimal parameter w^* .

Vectorized Version The *Graphics processing unit* (GPU) and linear algebra libraries of a modern computers make *vectorized* implementations much faster—i.e. writing your computations as vector or matrix products will make your code much faster than using for-loops. We can do this for the simple linear model above as follows. Firstly, regarding the data, we can write the collection of N features as $\mathbf{x} = [x_1, \dots, x_N]^T$, and similarly, the responses as $\mathbf{y} = [y_1, \dots, y_N]^T$. Now the vectorized form of the loss function is:

$$\ell(w; \mathcal{D}) = \frac{1}{N} \|\mathbf{w} \cdot \mathbf{x} - \mathbf{y}\|_2^2 \tag{5}$$

where $\|\cdot\|_2^2$ is the squared (Euclidean) two norm. Following the same derivation as above but keeping the vector notation, the optimal setting of the weights can then be written in vectorized form as: $w^* = (\mathbf{y}^T \mathbf{x}) / (\mathbf{x}^T \mathbf{x})$.

1.3 Maximum Likelihood Estimation: A General Recipe

While sensible, the above procedure we used for finding w^* could still seem arbitrary and unsound. For example, recalling that the goal of predictive modeling is to capture $\mathbb{P}(y|\mathbf{x})$, how does what we did relate to $\mathbb{P}(y|\mathbf{x})$? Moreover, are there other choices than the squared loss function? We will now give a general procedure for deriving optimization objectives known as *maximum likelihood estimation*.

Statistical Divergences Yet before introducing maximum likelihood estimation, we need to visit the concept of a statistical *divergence*. A divergence is like a loss function but applied to probability distributions. The most commonly employed divergence is the *Kullback-Leibler divergence* (KLD):

$$\text{KLD}[p(z)||q(z)] \triangleq \mathbb{E}_{p(z)} \left[\log \frac{p(z)}{q(z)} \right] = \int_z p(z) \left(\log \frac{p(z)}{q(z)} \right) dz,$$

where z is the random variable of interest, and we want to compare two distributions over z : $p(z)$ vs $q(z)$. The KLD is an information theoretic quantity that represents the number of bits lost when $q(z)$ is used to approximate $p(z)$. This means that the KLD is *not* symmetric: $\text{KLD}[p(z)||q(z)]$ does not necessarily equal $\text{KLD}[q(z)||p(z)]$, thus making the order of the arguments important. However, no matter the order of the arguments, the KLD will be exactly zero when $p(z) = q(z)$:

$$\text{KLD}[p(z)||p(z)] = \mathbb{E}_{p(z)} \left[\log \frac{p(z)}{p(z)} \right] = \mathbb{E}_{p(z)} [\log 1] = \log 1 = 0.$$

There are other divergences, such as the (squared) Hellinger divergence:

$$\mathcal{H}^2[p(z)||q(z)] \triangleq 1 - \int_z \sqrt{p(z) \cdot q(z)} dz.$$

The Hellinger divergence is symmetric, but unfortunately, it is less commonly employed due to it having an integral that is usually more difficult to evaluate. Both the Hellinger and KLD are members of the family of f -divergences.

Models as Probability Distributions Previously, we defined the model just as a generic function $f(\mathbf{x})$. Now we will be more particular interpreting $f(\mathbf{x})$, embedding it within a distribution function. This will, firstly, allow us to give a probabilistic interpretation to the model itself, unlocking operations such as marginalization, sampling, etc. Secondly, it will allow us to apply a statistical divergence to the model, directly quantifying the gap between the model and the true generative process $\mathbb{P}(y|\mathbf{x})$. To do this, we will pick a probability distribution $p(y;\theta)$, where y still denotes the label and θ denotes the parameter(s). For example, for the normal distribution $\theta = \{\mu, \sigma\}$, the mean and the standard deviation respectively. Since in our running example $y \in \mathbb{R}$, technically we can pick any distribution

with support over all the real numbers—e.g. normal, Laplace, student-t, etc.—but each comes with its own probabilistic assumptions. For example, the student-t distribution has heavier tails than the normal distribution, meaning that building a model with the student-t assumes that we will see more outlying points.

We will consider the general construction again in a later section. For now, let's assume that we choose $p(y; \theta)$ to be a normal distribution. Now taking our linear model $f(x) = w \cdot x$, we will use this model to parameterize just the mean μ :

$$\begin{aligned}
p(y; f(x)) &\triangleq \text{Normal}(y; \mu = f(x), \sigma) \\
&= \frac{1}{\sigma\sqrt{2\pi}} \exp\left\{-\frac{(\mu - y)^2}{2\sigma^2}\right\} \\
&= \frac{1}{\sigma\sqrt{2\pi}} \exp\left\{-\frac{(f(x) - y)^2}{2\sigma^2}\right\} \\
&= \frac{1}{\sigma\sqrt{2\pi}} \exp\left\{-\frac{(w \cdot x - y)^2}{2\sigma^2}\right\}.
\end{aligned} \tag{6}$$

The parameter σ will also have to be set somehow. We could set it with the same function, tying the mean and standard deviation, e.g. $\sigma = \exp\{f(x)\}$, where the $\exp\{\cdot\}$ ensures that the standard deviation is positive. However, this would mean that as the mean increases, so does the standard deviation, which is an assumption that would be inappropriate for many applications. Alternatively, we could set the standard deviation via a second linear model: $\sigma = \exp\{f'(x)\} = \exp\{u \cdot x\}$, where $u \in \mathbb{R}$ is another parameter that defines this second linear model. In the statistics literature, regression models that have a σ that is constant w.r.t. x are called *homoskedastic*. If σ varies with x , then the model is called *heteroskedastic*. Both of the cases above, where $\sigma = \exp\{f(x)\}$ or $\sigma = \exp\{f'(x)\}$, are heteroskedastic. In this course, for simplicity, we will often assume the models are homoskedastic.

Divergence as an Optimization Objective We now have the pieces in place to derive the maximum likelihood estimation procedure—the standard procedure we will use to train models throughout the course. Maximum likelihood estimation means that we seek to minimize the KLD between the true generative distribution and our probabilistic predictive model:

$$\begin{aligned}
\text{KLD} [\mathbb{P}(y|x) \parallel p(y; f(x))] &= \mathbb{E}_{\mathbb{P}(y|x)} \left[\log \frac{\mathbb{P}(y|x)}{p(y; f(x))} \right] \\
&= \underbrace{\mathbb{E}_{\mathbb{P}(y|x)} [\log \mathbb{P}(y|x)]}_{-\mathbb{H}[\mathbb{P}(y|x)]} - \mathbb{E}_{\mathbb{P}(y|x)} [\log p(y; f(x))] \\
&= \mathbb{E}_{\mathbb{P}(y|x)} [-\log p(y; f(x))] - \underbrace{\mathbb{H} [\mathbb{P}(y|x)]}_{\text{constant w.r.t. } p(y; f(x))}
\end{aligned} \tag{7}$$

where $\mathbb{H}[p(x)] = \int_x p(x)(-\log p(x))dx$ denotes the differential entropy of the distribution $p(x)$. $\mathbb{H} [\mathbb{P}(y|x)]$ denotes the entropy of the true generating process $\mathbb{P}(y|x)$, and thus it does not involve the model. This matters because we will eventually optimize this KLD w.r.t. $f(x)$, and in turn, terms that are not a function of $f(x)$ ‘fall out of’ the optimization problem. We can see this explicitly when we take the derivative w.r.t. the model parameters, since $\frac{d}{dw} \mathbb{H} [\mathbb{P}(y|x)] = 0$.

Yet, notice that Equation 7 involves a particular value of the features \mathbf{x} . Or in other words, both the model and true distribution are conditioned on a particular feature value, and we are evaluating the KLD only at those features. Of course, in the training data, we have multiple feature observations. This means that there is another distribution to be concerned with: the distribution that generates the features, $\mathbb{P}(\mathbf{x})$. We do not model this distribution directly. Rather, we will incorporate it into the maximum likelihood formulation by adding an outer expectation over $\mathbb{P}(\mathbf{x})$:

$$\mathbb{E}_{\mathbb{P}(\mathbf{x})} \text{KLD} [\mathbb{P}(y|\mathbf{x}) \parallel p(y; f(\mathbf{x}))] = \mathbb{E}_{\mathbb{P}(\mathbf{x})} \mathbb{E}_{\mathbb{P}(y|\mathbf{x})} [-\log p(y; f(\mathbf{x}))] - \mathbb{E}_{\mathbb{P}(\mathbf{x})} [\mathbb{H} [\mathbb{P}(y|\mathbf{x})]]$$

Putting everything together and taking the model parameter to again be \mathbf{w} , we have the complete maximum likelihood optimization problem:

$$\begin{aligned} w^* &= \arg \min_{\mathbf{w}} \mathbb{E}_{\mathbb{P}(\mathbf{x})} \text{KLD} [\mathbb{P}(y|\mathbf{x}) \parallel p(y; f(\mathbf{x}; \mathbf{w}))] \\ &= \arg \min_{\mathbf{w}} \mathbb{E}_{\mathbb{P}(\mathbf{x})} \mathbb{E}_{\mathbb{P}(y|\mathbf{x})} [-\log p(y; f(\mathbf{x}; \mathbf{w}))] - \mathbb{E}_{\mathbb{P}(\mathbf{x})} [\mathbb{H} [\mathbb{P}(y|\mathbf{x})]] \\ &= \arg \min_{\mathbf{w}} \mathbb{E}_{\mathbb{P}(\mathbf{x})} \mathbb{E}_{\mathbb{P}(y|\mathbf{x})} [-\log p(y; f(\mathbf{x}; \mathbf{w}))] \end{aligned} \quad (8)$$

where, again, the entropy term drops out because it does not involve the model and in turn, the parameter \mathbf{w} that we are optimizing.

Monte Carlo Approximation There is one remaining obstacle to solving the optimization problem in Equation 8. It requires taking the expectation w.r.t. the true generative process, $\mathbb{P}(y, \mathbf{x})$. As stated above, we do not have access to this distribution except through the samples that constitute our training data: $y \sim \mathbb{P}(y|\mathbf{x})$ and $\mathbf{x} \sim \mathbb{P}(\mathbf{x})$. Fortunately, this allows us to compute what's called a *Monte Carlo* (MC) approximation of the expectation; for a generic distribution $P(\mathbf{x})$ and a function of the random variable $\phi(\mathbf{x})$, this is:

$$\mathbb{E}_{P(\mathbf{x})} [\phi(\mathbf{x})] \approx \frac{1}{S} \sum_{s=1}^S \phi(x_s), \quad x_s \sim P(\mathbf{x}),$$

where $S \in \mathbb{N}^+$ is the number of samples. While this approximation is valid for any number of samples (above zero), the approximation becomes better and better as $S \rightarrow \infty$, becoming exact only asymptotically. Applying the MC expectation to Equation 8, we have:

$$\begin{aligned} w^* &= \arg \min_{\mathbf{w}} \mathbb{E}_{\mathbb{P}(\mathbf{x})} \mathbb{E}_{\mathbb{P}(y|\mathbf{x})} [-\log p(y; f(\mathbf{x}; \mathbf{w}))] \\ &\approx \arg \min_{\mathbf{w}} \frac{1}{N} \sum_{n=1}^N -\log p(y_n; f(x_n; \mathbf{w})) \end{aligned} \quad (9)$$

where the sum is over the training samples $\{(x_n, y_n)\}_{n=1}^N$. It follows that, the larger training set we have, the better we should be approximating our true optimization target, $\mathbb{E}_{\mathbb{P}(\mathbf{x})} \text{KLD} [\mathbb{P}(y|\mathbf{x}) \parallel p(y; f(\mathbf{x}; \mathbf{w}))]$.

Deriving the Squared Loss Function We will now work through an end-to-end derivation and eventually arrive at the same loss function used in Equation 2. Keeping with the same setup, we assume $f(\mathbf{x}; \mathbf{w}) = \mathbf{w} \cdot \mathbf{x}$ and $p(y; f(\mathbf{x}; \mathbf{w})) = \mathcal{N}(y; \mu = f(\mathbf{x}; \mathbf{w}), \sigma = 1)$, where

the normal distribution's variance is fixed at one (i.e. the homoskedastic assumption). The final optimization problem is then:

$$\begin{aligned}
w^* &= \arg \min_w \mathbb{E}_{\mathbb{P}(x)} \text{KLD} [\mathbb{P}(y|x) \parallel p(y; f(x; w))] \\
&= \arg \min_w \mathbb{E}_{\mathbb{P}(x)} \mathbb{E}_{\mathbb{P}(y|x)} [-\log p(y; f(x; w))] \quad (\text{drop entropy term}) \\
&\approx \arg \min_w \frac{1}{N} \sum_{n=1}^N -\log p(y_n; f(x_n; w)) \quad (\text{Monte Carlo approximation}) \\
&= \arg \min_w \frac{1}{N} \sum_{n=1}^N -\log \text{N}(y_n; \mu_n = f(x_n; w), \sigma = 1) \\
&= \arg \min_w \frac{1}{N} \sum_{n=1}^N -\log \left\{ \frac{1}{\sigma \sqrt{2\pi}} \exp \left\{ -\frac{(\mu_n - y_n)^2}{2\sigma^2} \right\} \right\} \\
&= \arg \min_w \frac{1}{N} \sum_{n=1}^N -\log \left\{ \frac{1}{\sqrt{2\pi}} \exp \left\{ -\frac{(f(x_n; w) - y_n)^2}{2} \right\} \right\} \tag{10} \\
&= \arg \min_w \frac{1}{N} \sum_{n=1}^N -\log \exp \left\{ -\frac{(f(x_n; w) - y_n)^2}{2} \right\} + \log \left\{ \sqrt{2\pi} \right\} \\
&= \arg \min_w \frac{1}{N} \sum_{n=1}^N \frac{1}{2} (f(x_n; w) - y_n)^2 + \log \left\{ \sqrt{2\pi} \right\} \\
&= \arg \min_w \frac{1}{2} \frac{1}{N} \sum_{n=1}^N (f(x_n; w) - y_n)^2 \quad (\text{drop } \sqrt{2\pi} \text{ constant}) \\
&= \arg \min_w \frac{1}{2} \frac{1}{N} \sum_{n=1}^N (w \cdot x_n - y_n)^2
\end{aligned}$$

where the $\log \left\{ \sqrt{2\pi} \right\}$ term is dropped because it does not depend on the parameters w . Thus, the maximum likelihood perspective gives us a more principled justification for use of the squared loss function as well as making its probabilistic assumptions explicit. If we were optimizing a heteroskedastic model w.r.t. the parameters controlling the variance, then this term would not drop from the optimization problem. The last line above is nearly the same as the final form of Equation 2 except for the constant $1/2$. Yet the presence of this constant does not change the solution, as the *maximum likelihood estimator* (MLE) for w^* is still the same as derived in Equation 4. In fact, the derivative becomes a bit simpler as the $1/2$ cancels the factor of 2 that is introduced when applying the power rule.

1.4 Generalized Linear Models

Above we discussed how a predictive model can be thought of as a probability distribution, with a specific function determining the mean variable. Specifically, we chose $p(y; f(x)) = \text{N}(y; \mu = f(x), \sigma)$. This flexible, modular framework allows us to construct models that meet our constraints or assumptions for the problem at hand. We call a linear model of the following form a *generalized linear model* (GLM):

$$\mathbb{E}_p[y|x] = g^{-1}(f(x; w)) = g^{-1}(w \cdot x) \tag{11}$$

where $f(x; w) = w \cdot x$, the linear model, and $g(\cdot)$ is known as the *link function*. In turn, $g^{-1}(\cdot)$ is called the *inverse link function*. Firstly, notice that in our running example of $p(y; f(x)) = N(y; \mu = f(x), \sigma)$, there is no g function since $\mu = f(x)$. Or to put it precisely, the link function is simply the identity function. This works in this case since the valid values of μ match the range of $f(x)$, namely all real numbers \mathbb{R} . But this will not be the case for all models of interest. Consider binary responses $y \in \{0, 1\}$. A common distribution with support over binary vales is the *Bernoulli* distribution: $p(y; \pi) = \pi^y \cdot (1 - \pi)^{1-y}$, where $\pi \in [0, 1]$ is the mean parameter. We cannot set $\pi = w \cdot x$ in this case since this would mean $\pi \in (-\infty, \infty)$, breaking the definition of the Bernoulli distribution. To fix this issue, we need to choose a g^{-1} function that transforms the output of $f(x)$ onto $(0, 1)$. We will examine one popular implementation for the Bernoulli case in the next section. If $y \in \mathbb{N}_{\geq 0}$, the non-negative natural numbers, the Poisson distribution is a commonly employed distribution with this support: $p(y; \lambda) = \lambda^y e^{-\lambda} / y!$, where $\lambda \in (0, \infty)$. In this case, it is common to choose the inverse link g^{-1} as the exponential function: $\lambda = \exp\{f(x)\}$.

Maximum Likelihood Derivative For the GLM under the maximum likelihood objective, we can write a general form for the chain rule derivative since it will always have the same four components:

$$\begin{aligned} \frac{d}{dw} \mathbb{E}_{\mathbb{P}(x)} \text{KLD} [\mathbb{P}(y|x) \parallel p(y; g^{-1} \circ f(x; w))] = \\ \left(\frac{d}{dp} \mathbb{E}_{\mathbb{P}(x)} \text{KLD} [\mathbb{P}(y|x) \parallel p(y)] \right) \left(\frac{d}{dg^{-1}} p(y; g^{-1}) \right) \left(\frac{d}{df} g^{-1}(f) \right) \left(\frac{d}{dw} f(x; w) \right). \end{aligned}$$

While for linear models $\frac{d}{dw} f(x; w)$ is usually easy to evaluate, this will be the most computationally intensive term for deep learning models.

1.5 Univariate Logistic Regression for Binary Labels

We will now examine a particular GLM mentioned in the preceding section—namely, a model for binary labels known as *logistic regression*. Assuming $y \in \{0, 1\}$, we can define the following predictive model:

$$\begin{aligned} p(y; f(x; w)) &= \text{Bernoulli}(y; \pi = g^{-1}(f(x; w))) \\ &= \pi^y \cdot (1 - \pi)^{1-y} \\ &= g^{-1}(f(x; w))^y \cdot (1 - g^{-1}(f(x; w)))^{1-y} \\ &= g^{-1}(w \cdot x)^y \cdot (1 - g^{-1}(w \cdot x))^{1-y}. \end{aligned} \tag{12}$$

Now we need to select the form of g^{-1} such that $g^{-1} : \mathbb{R} \mapsto (0, 1)$. The most common choice of g^{-1} is the *logistic function*, which is where the name *logistic regression* originates:

$$\text{logistic}(z) \triangleq \frac{1}{1 + \exp\{-z\}}.$$

We call this the *logistic function* because it is the cumulative distribution function of the standard logistic distribution (‘standard’ meaning location 0, scale 1). Plugging this into

the above expressions gives us our final form of logistic regression:

$$\begin{aligned}
p(y; f(x; w)) &= \text{Bernoulli}(y; \pi = \text{logistic}(f(x; w))) \\
&= \text{logistic}(f(x; w))^y \cdot (1 - \text{logistic}(f(x; w)))^{1-y} \\
&= \left(\frac{1}{1 + \exp\{-f(x; w)\}} \right)^y \cdot \left(1 - \frac{1}{1 + \exp\{-f(x; w)\}} \right)^{1-y} \\
&= \left(\frac{1}{1 + \exp\{-w \cdot x\}} \right)^y \cdot \left(1 - \frac{1}{1 + \exp\{-w \cdot x\}} \right)^{1-y}.
\end{aligned} \tag{13}$$

Plugging this model into the maximum likelihood objective, we have:

$$\begin{aligned}
w^* &= \arg \min_w \mathbb{E}_{\mathbb{P}(x)} \text{KLD}[\mathbb{P}(y|x) \parallel \text{Bernoulli}(y; \pi = \text{logistic}(f(x; w)))] \\
&= \arg \min_w \mathbb{E}_{\mathbb{P}(x)} \mathbb{E}_{\mathbb{P}(y|x)} [-\log \text{Bernoulli}(y; \pi = \text{logistic}(f(x; w)))] \quad (\text{drop entropy term}) \\
&\approx \arg \min_w \frac{1}{N} \sum_{n=1}^N -\log \text{Bernoulli}(y_n; \pi = \text{logistic}(f(x_n; w))) \quad (\text{Monte Carlo approximation}) \\
&= \arg \min_w \frac{1}{N} \sum_{n=1}^N -\log \{ \text{logistic}(f(x_n; w))^{y_n} \cdot (1 - \text{logistic}(f(x_n; w)))^{1-y_n} \} \\
&= \arg \min_w \frac{1}{N} \sum_{n=1}^N -y_n \log \text{logistic}(f(x_n; w)) - (1 - y_n) \log (1 - \text{logistic}(f(x_n; w))) \\
&= \arg \min_w \frac{1}{N} \sum_{n=1}^N -y_n \log \text{logistic}(w \cdot x_n) - (1 - y_n) \log (1 - \text{logistic}(w \cdot x_n)) \\
&= \arg \min_w \frac{1}{N} \sum_{n=1}^N -y_n \log \left(\frac{1}{1 + \exp\{-w \cdot x\}} \right) - (1 - y_n) \log \left(1 - \left(\frac{1}{1 + \exp\{-w \cdot x\}} \right) \right).
\end{aligned}$$

This does not have as intuitive a form as the squared error loss function we examined earlier, but this expression is known as the *binary cross-entropy loss*. The ‘cross-entropy’ part is a bit of mis-characterization of this particular loss. Cross-entropy is defined as $\mathbb{H}[p(x), q(x)] = -\mathbb{E}_{p(x)} [\log q(x)]$ for two distributions $p(x)$ and $q(x)$, which is what we have in step #2 of the above derivation but also in step #2 of Equation 10. Thus, squared error could also be called a ‘cross-entropy loss’, but when people say that, they usually are assuming the labels take on binary or (as we’ll see later) categorical values.

1.6 Gradient Descent

If you try to take the derivative of the cross-entropy loss above, set it to zero, and solve for w , you will find that you cannot isolate w to one side of the equation, meaning that the optimization problem has no ‘closed-form’ solution. Instead we need to find a numerical solution that will only approximate the true optimum. We will use a procedure known as ‘gradient descent’, a.k.a. ‘steepest descent’. The intuition is that we’ll start with an initial guess at the value of w and slowly walk down the loss surface, following the direction of steepest descent according to the derivative at our current point. For a generic function

$\phi(z)$ that we wish to minimize, we can apply gradient descent by iterating the equation:

$$z_{t+1} = z_t - \alpha \cdot \frac{d}{dz_t} \phi(z_t)$$

where $\alpha > 0$ is the *learning rate* or *step size* that controls how aggressively we move down the path of steepest descent at each iteration. It is common to set α to be large for the early iterations and then slowly decay it, taking smaller and smaller step sizes, when the procedure gets close to a minimum. We can know if we've approximately converged by tracking the derivative $\frac{d}{dz_t} \phi(z_t)$ since it should be exactly zero at the minimum (or any other critical point—a topic for later in the course).

Applying Gradient Descent to Logistic Regression We will now apply gradient descent to the logistic regression model. The first step will be to take the derivative of the cross-entropy loss w.r.t. the parameter w . When doing this, we will use the fact that the derivative of the logistic function is: $d/dz \text{logistic}(z) = \text{logistic}(z) \cdot (1 - \text{logistic}(z))$. I'll leave showing this to the reader as an exercise.

$$\begin{aligned} \frac{d}{dw} \ell(w; \mathcal{D}) &= \frac{d}{dw} \left[\frac{1}{N} \sum_{n=1}^N -\log \text{Bernoulli}(y_n; \pi = \text{logistic}(f(x_n; w))) \right] \\ &= \frac{1}{N} \sum_{n=1}^N -y_n \frac{d}{dw} [\log \text{logistic}(w \cdot x_n)] - (1 - y_n) \frac{d}{dw} [\log (1 - \text{logistic}(w \cdot x_n))] \\ &= \frac{1}{N} \sum_{n=1}^N -y_n \cdot \frac{\text{logistic}(w \cdot x_n) \cdot (1 - \text{logistic}(w \cdot x_n))}{\text{logistic}(w \cdot x_n)} \cdot x_n \\ &\quad - (1 - y_n) \frac{-\text{logistic}(w \cdot x_n) \cdot (1 - \text{logistic}(w \cdot x_n))}{1 - \text{logistic}(w \cdot x_n)} \cdot x_n \\ &= \frac{1}{N} \sum_{n=1}^N -y_n \cdot (1 - \text{logistic}(w \cdot x_n)) \cdot x_n + (1 - y_n) \text{logistic}(w \cdot x_n) \cdot x_n \\ &= \frac{1}{N} \sum_{n=1}^N (\text{logistic}(w \cdot x_n) - y_n) \cdot x_n \\ &= \frac{1}{N} \sum_{n=1}^N (\mathbb{E}_p[y|x_n] - y_n) \cdot x_n \quad (\text{via definition of the logistic regression GLM}). \end{aligned}$$

Using the definition of logistic regression as a GLM reveals that the derivative takes on a very sensible form: the difference between the expected value of y and its actual value, y_n , scaled by the feature value x_n . When $(\mathbb{E}_p[y|x_n] - y_n) \approx 0$, the model is a *confident* and *accurate* predictor and thus the contribution of this data point is negligible. If the total derivative is zero, then the model is accurately predicting all training points with *maximal* confidence. Plugging this derivative into the gradient descent equation, we have:

$$\begin{aligned} w_{t+1} &= w_t - \alpha \cdot \frac{d}{dw_t} \ell(w_t; \mathcal{D}) \\ &= w_t - \alpha \left[\frac{1}{N} \sum_{n=1}^N (\text{logistic}(w_t \cdot x_n) - y_n) \cdot x_n \right]. \end{aligned} \tag{14}$$

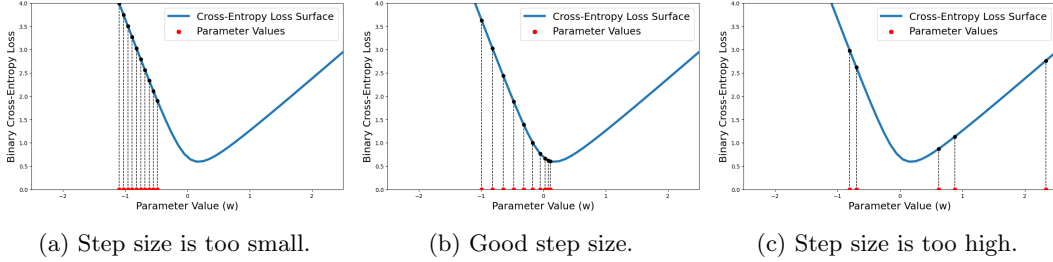


Figure 1: *Gradient descent for logistic regression, with varying step sizes.*

For each run, gradient descent is run for 10 steps with a fixed step size. In (a), the step size is too small ($\alpha = .02$); in (b), the step size is suitable ($\alpha = .05$); in (c), the step size is too large ($\alpha = 1.5$).

We would implement the above equation numerically by guessing an initial value, w_0 , making a prediction using w_0 (i.e. evaluate the logistic output) for all training points, compute the derivative by combining the predictions, the labels and features, as shown above, and lastly updating the weight to arrive at w_1 . That process is then repeated for a maximum number of rounds or until the derivative is sufficiently close to zero (e.g. 1×10^{-4}).

Figure 1 shows a simulation for a one-parameter logistic regression model. The (binary) cross-entropy loss surface is visualized by the blue line. The intermediate parameter estimates (w_t) are shown along the x-axis in red, and the dotted line connects them to the loss value that they resulted in. Subfigure 1a shows when the step size is too small, as the maximum number of iterations is reached before the minimum is found. Subfigure 1 shows the other extreme: the step size is too big and so the minimum cannot be found. Instead, the optimizer ‘jumps’ over the minimum and flies off to the right-hand side of the plot. Subfigure 1b shows when the step size is properly set, as the optimizer gracefully descends to the minimum.

1.7 Models of Artificial Neurons & the Perceptron

So far, the narrative of these notes takes a purely statistical perspective. However, deep learning also has roots in artificial intelligence, which at times has taken inspiration from biological intelligence. The first major step towards formulating a computational model of biological neurons was taken in 1943 by McCulloch and Pitts. A biological neuron can be, very roughly, thought of as a model that takes an input signal (dendrite), performs computation (soma), and passes the output to other connected neurons (axon to synapse to other neuron’s dendrite). A visualization is shown in Figure 2a. McCulloch and Pitts proposed a model whose input can be either *excitatory* or *inhibitory*. Its output can then be either *quiet* or *firing*, depending on if the the number of excitatory inputs is equal to or greater than some fixed threshold.

In 1957, Frank Rosenblatt famously implemented a very similar model on an IBM 704 computer, calling it *the perceptron*, and demonstrated that it could ‘learn’ to classify simple patterns. This demonstration received wide media attention; the New York Times had an article with the headline: “Electronic ‘Brain’ Teaches Itself.” The perceptron can be defined

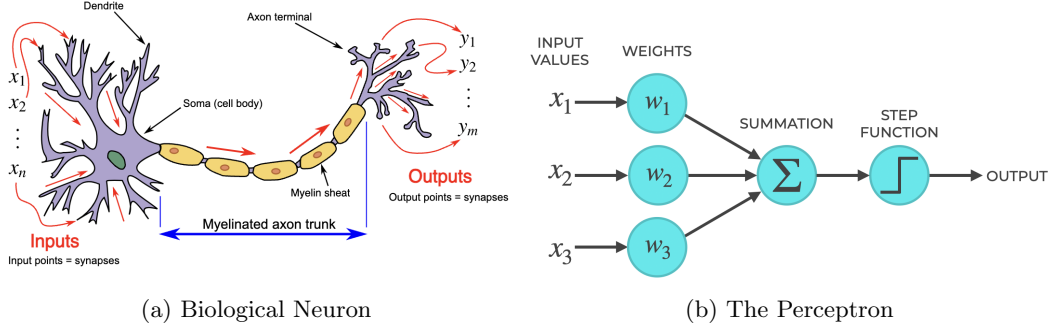


Figure 2: *Biological vs Artificial Neurons.*

for an input $x \in \mathbb{R}$ and parameters $w \in \mathbb{R}$ and $b \in \mathbb{R}$ as:

$$\hat{y} = \psi(w \cdot x + b), \text{ where } \psi(z) = -1 + 2 \cdot \mathbb{I}[z > 0],$$

with $\mathbb{I}[\cdot]$ denoting the indicator function that evaluates to 1 if its argument is true and 0 otherwise. In turn, $\psi(z) = +1$ if its input is greater than zero and -1 otherwise. Thus \hat{y} is the prediction, and it should model binary data represented as $y \in \{-1, +1\}$. The similarity to the McCulloch and Pitts model is that the input x is modulated by a weight w , and if $w \cdot x > -b$, then \hat{y} ‘fires,’ having a value of one. Thus $-b$ serves as the activation threshold proposed by McCulloch and Pitts; though it was discrete in the original model and real-valued here. A visualization of the perceptron is shown in Figure 2b.

The perceptron learning algorithm is as follows for an N -sized data set $\mathcal{D} = \{x_n, y_n\}_{n=1}^N$, $x \in \mathbb{R}$, $y \in \{-1, +1\}$. The algorithm starts by randomly initializing w and b ; denote these values as w_0 and b_0 . Then for every time step t , a training pair is selected at random (x_r, y_r) , the perceptron’s prediction is computed as $\hat{y}_r = \psi(w_0 \cdot x_r + b_0)$, and this prediction is checked for correctness, $\hat{y}_r = y_r$. If the prediction is correct, the process repeats for another feature-label pair. If the prediction is incorrect, then the following update is performed to obtain w_{t+1} and b_{t+1} :

$$w_{t+1} = w_t + y_r \cdot x_r, \quad b_{t+1} = b_t + y_r.$$

This process is repeated for either a maximum number of iterations or until all points in the training data have their label correctly predicted.

The update rule for the perceptron learning algorithm looks a bit mysterious, and it is certainly not clear, at least at a first glance, why iterating them amounts to ‘learning.’ Yet we can examine some parallels to the steepest descent equations for logistic regression to gain intuition. Let’s consider the logistic regression update rule given in Section 1.6. But we will make two simplifying assumptions: the update will be computed for just one data point and with a step size of $\alpha = 1$:

$$w_{t+1} = w_t - (\mathbb{E}[y|x] - y) \cdot x = w_t + (y - \mathbb{E}[y|x]) \cdot x.$$

By comparing this equation to the perceptron’s weight update, the only difference is the factor that is multiplied with the features x : $(y - \mathbb{E}[y|x])$ for logistic regression (for $y \in \{0, 1\}$) vs y for the perceptron. Now let’s assume logistic regression’s output can only take

on the extreme values of 0 and 1, i.e. $\mathbb{E}[y|x] \in \{0, 1\}$. Thus when $\mathbb{E}[y|x] = y$, their difference will be zero and there will be no update to the parameters. When $\mathbb{E}[y|x] \neq y$, the difference will either be +1 when $y = 1$ or -1 when $y = 0$. Notice that the perceptron learning rule is then recovered *exactly* due to the label support being defined as $\{-1, +1\}$. The update for the second parameter, b , can be recovered by thinking of the features as being a two dimensional vector: $\mathbf{x} = [x_0, x_1]^T$ with x_0 always being set to one. Then given the parameter vector $\boldsymbol{\theta} = [b, w]^T$, we have $\boldsymbol{\theta}^T \mathbf{w} = b \cdot x_0 + w \cdot x_1 = w \cdot x_1 + b$. We will discuss this representation in more detail in the next section. In general, we can think of the perceptron as performing a ‘hard’ update since its predictions are either completely correct or incorrect. On the other hand, logistic regression’s learning rule is more precise since it has a more granular notion of how far off its predictions (as represented by $\mathbb{E}[y|x]$) are from the true label.

While the invention of the perceptron generated much excitement, this excitement was short-lived. In 1969, Marvin Minsky and Seymour Papert published a book entitled, *Perceptrons: An Introduction to Computational Geometry*. The authors demonstrated the perceptron’s limitations—namely, that it could only learn linear decision functions. Thus the perceptron could not model a simple logical function like XOR (i.e. *exclusive or*). This dissolved enthusiasm for artificial neurons as a potential path towards artificial intelligence, and symbolic approaches began to be favored instead. This period generally saw the decline in funding for artificial intelligence in the 1970’s and 1980’s. When similar approaches started to gain traction again in the 1990’s, they were often re-branded as ‘machine learning’ instead of ‘artificial intelligence.’

2 Supervised Learning with Multiple Linear Regression

In the previous section, we have built up our foundation for specifying statistical predictive models and their loss functions. We also covered how to train them using steepest descent. Now we will look to expand these models such that they can have multiple input features and multiple output features. The underlying principles will primarily be the same, but now as parameters and outputs will be vector-valued, we will need to use vector calculus for deriving the corresponding learning updates.

2.1 Multiple Features and Feature Expansions

We will now consider regression models for which the input features are multi-dimensional. In turn, we will now need to specify multiple weight parameters—one for each input dimension. Thus, our (generalized) regression models will be specified as:

$$\mathbb{E}[y|\mathbf{x}] \triangleq g^{-1}(\mathbf{w}^T \mathbf{x}), \text{ where } \mathbf{w} = \begin{bmatrix} w_1 \\ \vdots \\ w_D \end{bmatrix} \in \mathbb{R}^D \text{ and } \mathbf{x} = \begin{bmatrix} x_1 \\ \vdots \\ x_D \end{bmatrix} \in \mathbb{R}^D, \quad (15)$$

and where $g^{-1}(\cdot)$ is still the (inverse) link function and y is the scalar label. The most obvious case where this formulation is applicable is for applications that have multiple feature observations. For example, y might be a binary indicator representing the presence

of heart disease, and \mathbf{x} could be various health features such as age, blood pressure, weight, cortisol levels, etc. We will now consider two other examples.

Vector Representation of Offset Parameter One simple but common use for the vector representation, even when there is only input feature, is to encode the offset a.k.a. bias parameter. In the perceptron model, this was b . The regression models we considered so far (except the perceptron) had just one parameter, meaning that we can interpret it as modeling the slope of a line. However, we might also want to model the bias / offset from the x -axis. Or in logistic regression, the offset determines where the decision boundary of $p(y|\mathbf{x}) = 0.5$ resides. We can treat the offset parameter like any other weight by appending a constant value of one to every feature vector. Thus we have:

$$\mathbb{E}[y|\mathbf{x}] \triangleq g^{-1}(w_1 \cdot x + w_0) = g^{-1}(\mathbf{w}^T \mathbf{x}), \text{ where } \mathbf{w} = \begin{bmatrix} w_0 \\ w_1 \end{bmatrix} \in \mathbb{R}^D \text{ and } \mathbf{x} = \begin{bmatrix} 1 \\ x \end{bmatrix} \in \mathbb{R}^D,$$

where w_0 is the offset parameter and w_1 is the weight that represents the slope parameter, like before. While this may seem to be just a notational trick, this allows for cleaner programmatic implementations since the offset parameter does not have to be treated differently than the other weight parameters. In this course, we will usually assume offset parameters are treated in this way, by appending a constant 1 to the feature vectors.

Polynomial Basis Expansion Linear models may seem limited, at first glance, to encoding only straight lines. Yet, one trick to make linear models more expressive while not adding much computational overhead is to use a *feature expansion*. Let's assume we have a scalar data $x \in \mathbb{R}$ and $y \in \mathbb{R}$. Our previous models for this data used only one parameter and thus encoded the slope of a line through the origin. Yet what if the relationship between the features and labels is more complex than that? One simple trick is to create 'dummy' features by replicating the existing feature. One simple but still very expressive way is to use a K -degree *polynomial basis*. This means that we create a new feature vector by taking x to powers up to degree K :

$$\tilde{\mathbf{x}} = [x^0 \ x^1 \ x^2 \ \dots \ x^K]^T = [1 \ x \ x^2 \ \dots \ x^K]^T.$$

Then when we use this expanded feature set in the linear model, the model encodes a K -degree polynomial, with the $K + 1$ parameters / weights serving as the coefficients:

$$\mathbb{E}[y|\mathbf{x}] \triangleq g^{-1}(\mathbf{w}^T \tilde{\mathbf{x}}) = g^{-1}\left(\sum_{k=0}^K w_k \cdot x^k\right), \text{ where } \mathbf{w} = \begin{bmatrix} w_0 \\ \vdots \\ w_K \end{bmatrix} \in \mathbb{R}^{K+1} \text{ and } \tilde{\mathbf{x}} = \begin{bmatrix} 1 \\ \vdots \\ x^K \end{bmatrix} \in \mathbb{R}^{K+1}.$$

Figure 3 shows a simulation in which the true function we are trying to model is a $K = 7$ degree polynomial. The plots show when real-valued regression is run with degrees ranging from one (linear) to fifteen. We see that once the degree of the model surpasses the degree of the true function ($K = 10$), then the model is well-fit to the data.

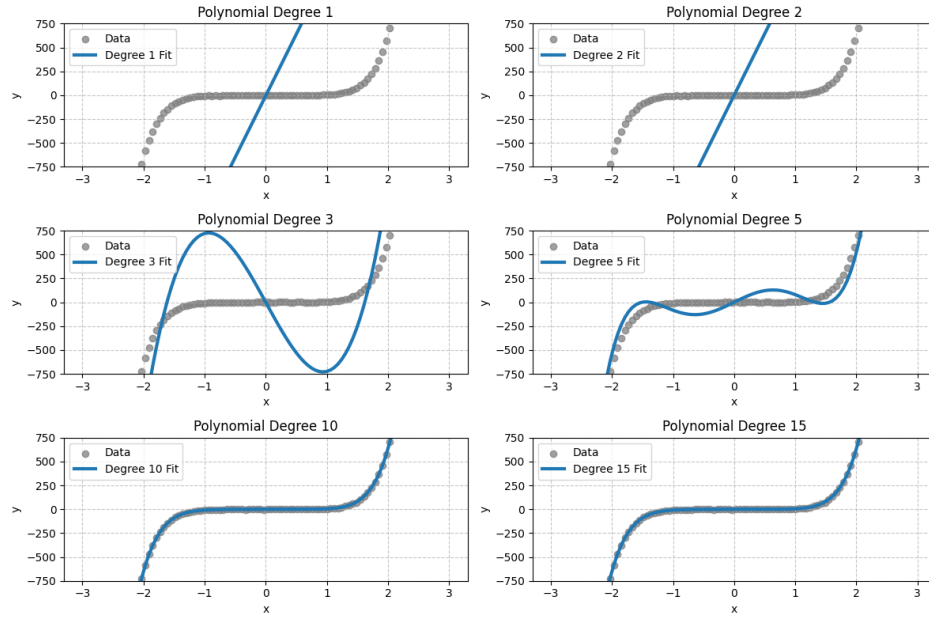


Figure 3: *Linear Regression with Polynomial Basis*. The data represents a 7th-degree polynomial. We see that as the degree of the model's feature basis increases, the model can express richer and richer functions.

2.2 Revisiting Gradient Descent with Multivariate Derivatives

2.3 Multiple Output Dimensions

3 Model Evaluation, Model Selection, and Regularization

4 Feedforward Neural Networks

4.1 Adaptive Basis Functions

5 Convolutional Neural Networks

6 Recurrent Neural Networks

7 Transformers

8 Deep Generative Models

9 Uncertainty Estimation

Bibliography