

are different. In software, we often phrase both as minimizing a cost function. Maximum likelihood thus becomes minimization of the negative log-likelihood (NLL), or equivalently, minimization of the cross entropy. The perspective of maximum likelihood as minimum KL divergence becomes helpful in this case because the KL divergence has a known minimum value of zero. The negative log-likelihood can actually become negative when  $\mathbf{x}$  is real-valued.

### 5.5.1 Conditional Log-Likelihood and Mean Squared Error

The maximum likelihood estimator can readily be generalized to the case where our goal is to estimate a conditional probability  $P(\mathbf{y} \mid \mathbf{x}; \boldsymbol{\theta})$  in order to predict  $\mathbf{y}$  given  $\mathbf{x}$ . This is actually the most common situation because it forms the basis for most supervised learning. If  $\mathbf{X}$  represents all our inputs and  $\mathbf{Y}$  all our observed targets, then the conditional maximum likelihood estimator is

$$\boldsymbol{\theta}_{\text{ML}} = \arg \max_{\boldsymbol{\theta}} P(\mathbf{Y} \mid \mathbf{X}; \boldsymbol{\theta}). \quad (5.62)$$

If the examples are assumed to be i.i.d., then this can be decomposed into

$$\boldsymbol{\theta}_{\text{ML}} = \arg \max_{\boldsymbol{\theta}} \sum_{i=1}^m \log P(\mathbf{y}^{(i)} \mid \mathbf{x}^{(i)}; \boldsymbol{\theta}). \quad (5.63)$$

**Example: Linear Regression as Maximum Likelihood** Linear regression, introduced earlier in section 5.1.4, may be justified as a maximum likelihood procedure. Previously, we motivated linear regression as an algorithm that learns to take an input  $\mathbf{x}$  and produce an output value  $\hat{y}$ . The mapping from  $\mathbf{x}$  to  $\hat{y}$  is chosen to minimize mean squared error, a criterion that we introduced more or less arbitrarily. We now revisit linear regression from the point of view of maximum likelihood estimation. Instead of producing a single prediction  $\hat{y}$ , we now think of the model as producing a conditional distribution  $p(y \mid \mathbf{x})$ . We can imagine that with an infinitely large training set, we might see several training examples with the same input value  $\mathbf{x}$  but different values of  $y$ . The goal of the learning algorithm is now to fit the distribution  $p(y \mid \mathbf{x})$  to all of those different  $y$  values that are all compatible with  $\mathbf{x}$ . To derive the same linear regression algorithm we obtained before, we define  $p(y \mid \mathbf{x}) = \mathcal{N}(y; \hat{y}(\mathbf{x}; \mathbf{w}), \sigma^2)$ . The function  $\hat{y}(\mathbf{x}; \mathbf{w})$  gives the prediction of the mean of the Gaussian. In this example, we assume that the variance is fixed to some constant  $\sigma^2$  chosen by the user. We will see that this choice of the functional form of  $p(y \mid \mathbf{x})$  causes the maximum likelihood estimation procedure to yield the same learning algorithm as we developed before. Since the