

Machine learning theory

Vilde Flusgrud

April 28, 2018

1 Introduction

- Estimation (try to determine a model(?) parameter) vs. prediction (try to predict the value of a random variable)
- The review is organized as follows. We begin by introducing polynomial regression as a simple example that highlights many of the core ideas of ML. The next few chapters introduce the language and major concepts needed to make these ideas more precise including tools from statistical learning theory such as overfitting, the bias-variance tradeoff, regularization, and the basics of Bayesian inference. The next chapter builds on these examples to discuss stochastic gradient descent and its generalizations. We then apply these concepts to linear and logistic regression, followed by a detour to discuss how we can combine multiple statistical techniques to improve supervised learning, introducing bagging, boosting, random forests, and XG Boost. These ideas, though fairly technical, lie at the root of many of the advances in ML over the last decade. The review continues with a thorough discussion of supervised deep learning and neural networks, as well as convolutional nets. We then turn our focus to unsupervised learning. We start with data visualization and dimensionality reduction before proceeding to a detailed treatment of clustering. Our discussion of clustering naturally leads to an examination of variational methods and their close relationship with mean-field theory. The review continues with a discussion of deep unsupervised learning, focusing on energy-based models, such as Restricted Boltzmann Machines (RBMs) and Deep Boltzmann Machines (DBMs). Then we discuss two new and extremely popular modeling frameworks for unsupervised learning, generative adversarial networks (GANs) and variational autoencoders (VAEs). We conclude the review with an outlook and discussion of promising research directions at the intersection physics and ML.

2 Why is machine learning difficult?

- Ingredients:
 - \mathbf{X} = the dataset
 - $g(\mathbf{w})$ = the model = a function of the parameters \mathbf{w}
 - $\mathcal{C}(\mathbf{X}, g(\mathbf{w}))$ = the cost function, allows us to judge how well the model performs on the observations.
 - \mathbf{X}_{train} = 90% of \mathbf{X}
 - \mathbf{X}_{test} = 10% of \mathbf{X}
 - $\hat{\mathbf{w}} = \operatorname{argmin}_{\mathbf{w}} \{\mathcal{C}(\mathbf{X}_{train}, g(\mathbf{w}))\}$ (the model is fit by minimizing the cost function)
 - $E_{in} = \mathcal{C}(\mathbf{X}_{train}, g(\hat{\mathbf{w}}))$ (the in-sample error)
 - $E_{out} = \mathcal{C}(\mathbf{X}_{test}, g(\hat{\mathbf{w}}))$ (The out-of-sample error. Used to evaluate the performance of the model.)
 - We almost always find that $E_{in} \geq E_{out}$
- It may be at first surprising that the model that has the lowest out-of-sample error E_{out} usually *does not* have the lowest in-sample error E_{in} . At first glance, the observation that the model providing the best explanation for the current dataset probably will not provide the best explanation for future datasets is very counter-intuitive.

Moreover, the discrepancy between E_{in} and E_{out} becomes more and more important, as the complexity of our data, and the models we use to make predictions, grows. As the number of parameters in the model increases, we are forced to work in high-dimensional spaces. The “curse of dimensionality” ensures that many phenomena that are absent or rare in low-dimensional spaces become generic. For example, the nature of distance changes in high dimensions, as evidenced in the derivation of the Maxwell distribution in statistical physics where the fact that all the volume of a d -dimensional sphere of radius r is contained in a small spherical shell around r is exploited. Almost all critical points of a function (i.e., the points where all derivatives vanish) are saddles rather than maxima or minima (an observation first made in physics in the context of the p -spin spherical spin glass). For all these reasons, it turns out that for complicated models studied in ML, predicting and fitting are very different things.

- Polynomial regression
 - Probabilistic process that assigns a label y_i to an observation x_i . Data generated by drawing samples from the equation $y_i = f(x_i) + \eta_i$

- $f(x_i)$ = some fixed (but possibly unknown) function. "Function used to generate the data".
- η_i = a Gaussian, uncorrelated noise variable, such that
 - * $\langle \eta_i \rangle = 0$
 - * $\langle \eta_i \eta_j \rangle = \delta_{ij} \sigma^2$
- σ = the noise strength. ($\sigma = 0$ = noiseless case.)
- $g_\alpha(x; \mathbf{w}_\alpha)$ = a family of functions which depend on some parameters \mathbf{w}_α . These functions represent the *model class* that we are using to model the data and make predictions. We chose the model class without knowing $f(x)$. The model class encode the *features* we choose to represent the data. For polynomial regression we will consider three different model classes:
 - i $g_1(x; \mathbf{w}_1)$ = all polynomials of order 1 (two parameters)
 - ii $g_3(x; \mathbf{w}_3)$ = all polynomials up to order 3 (four parameters)
 - iii $g_{10}(x; \mathbf{w}_{10})$ = all polynomials of order 10 (eleven parameters)
- The different number of parameters reflects that the three models have different *model complexities*. Thinking of each term in the polynomial as a "feature", increasing the polynomial order increases the number of features.
- Fit the models on the generated training samples using standard least-squares regression.
- Observe that at small sample sizes, noise can create fluctuations in the data that look like genuine patterns. While simple models are forced to ignore them and focus on the larger trends, complex models can capture both the global trends and noise-generated patterns at the same time. The model can then be tricked into thinking the noise encodes real information = "**overfitting**" = a steep drop-off in predictive performance. Can guard against overfitting in two ways:
 - * Use less expressive models with fewer parameters
 - * Collect more data so the likelihood that the noise appears patterned decreases.

This relates to the **bias-variance** tradeoff: when training data limited, one can often get better predictive performance by using less expressive model rather than a complex one. The simpler model has more "bias" but is less dependent on the particular realization of the training data., i. e. less "variance". Some universal lessons:

- * Fitting is not predicting. Fitting existing data well is fundamentally different from making predictions about new data. (fitting = estimation? or not?)

- * Using a complex model can result in overfitting. Increasing a model's complexity will usually yield better results on the training data. However when the training data size is small and the data are noisy, this results in *overfitting* and can substantially degrade the predictive performance of the model.
- * For complex datasets and small training sets, simple models can be better at prediction than complex models due to the bias-variance tradeoff. It takes less data to train a simple model than a complex one. Therefore, even though the correct model is guaranteed to have better predictive performance for an infinite amount of training data (less bias), the training errors stemming from finite-size sampling (variance) can cause simpler models to out-perform the more complex model when sampling is limited.
- * It is difficult to generalize beyond the situations encountered in the training data.

3 Basics of statistical learning theory

- Goal: the sense in which learning is possible, with focus on supervised learning. Ingredients:
 - $y = f(x)$ = an unknown function
 - \mathcal{H} = a hypothesis set that we fix, consisting of all functions we are willing to consider, defined also on the domain of f . The set may be uncountably infinite (e.g. if there are real-valued parameters to fit). Our choices here depends usually on our intuition about the problem.
 - $(x_i, y_i), \quad i = 1 \dots N$ = a set of pairs produced by $f(x)$ which serve as observable data.
 - Our goal: find a function $h \in \mathcal{H}$ approximating $f(x)$ as best as possible, $h \approx f$ in some strict mathematical sense specified below. Then say we *learned* $f(x)$.
 - If $f(x)$ can in principle take any value on *unobserved* inputs, how is it possible to learn in any meaningful sense? Learning is possible in the restricted sense that the fitted model will probably perform approximately as well on new data as it did on training data.
 - E = appropriately chosen error function (e.g. sum of squared errors in linear regression)
 - When we are training we only have access to E_{in} (fitting). Our goal is to minimize E_{out} (predicting) - the performance on new data.

- Can we say something about the relationship between E_{in} and E_{out} ? Yes: it's the domain of statistical learning theory.

- Three schematics

- Figure 4: Shows E_{in} and E_{out} as functions of the amount of training data. Assumes large data amount and that the model cannot exactly fit the true function $f(x)$. In the infinite data limit, the two errors must approach the same value, which is our model's **bias** = the best our model could do given infinite training data to beat down sampling noise = a property of the kind of functions/model class we use to approximate $f(x)$. In general: more complex model class = smaller bias. But, do not have infinite data. Thus, better minimize E_{out} than the bias. Can decompose E_{out} into

- * The bias
- * The variance = measures the typical errors introduced in training our model due to sampling noise from having a finite training set.

Final quantity shown is the difference between E_{out} (generalization) and E_{in} (training). Measures difference between fitting and predicting. Models with large difference **overfit** the data. Statistical learning lesson: not enough to minimize E_{in} since E_{out} may still be large. This insight leads to the idea of **regularization**.

- Figure 5: Shows E_{out} as a function of model complexity = number of parameters/features f_{ex} , but not always = model complexity is a subtle idea, defining it precisely is one of the great achievements of statistical learning theory = roughly speaking, it is a measure of the complexity of the model class used to approximate $f(x)$. Considering a training data set of fixed size, E_{out} will be a non-monotonic function of the model complexity and generally minimized for models of *intermediate* complexity. Because: while using a more complex model always reduces the bias, at some point the model becomes too complex for the amount of training data and the generalization error becomes large due to high variance = may be more suitable to use a more biased model with small variance than a less-biased model with large variance = the **bias-variance tradeoff**.
- Figure 6: Another way to visualize the bias-variance tradeoff. Shows how a complex (high variance, low-bias) model vs a simpler (low variance, high bias) model lands compared to the true model.
 - * The complex exhibits larger fluctuations while its average will be closer to the true model. The simpler fluctuates less, but is

on average further from the true model. (This reminds one of the accuracy (achieved by the complex model) and precision (achieved by the simpler model) discussion in experimental physics).

- * In general, more complex model needs more training data. This is the cause of the larger fluctuations for the complex model. However, when increasing the data the more complex eventually performs better. Thus, the choice of complexity depends on the amount of training data.

- Bias-Variance Decomposition: dig further into the central principle of the bias-variance tradeoff. Expressiveness vs sensitivity to training data fluctuations. Oftentimes in physics, we are mostly concerned with expressivity, *e.g. whether the true ground state wave function can be well approximated by a class of variational wavefunctions such as a matrix product state*. In the learning context, there is the additional challenge of finding the best variational state with finite sampling. We will see that while this concept is a generally useful heuristic (technique) to keep in mind, it is a mathematically precise statement when decomposing the squared error. (what is meant by this entire sentence?) Finally, we note that a better term would be the bias-variance *decomposition*, as it is possible to have high bias *and* high variance. We'll discuss the b-v tradeoff in the continuous predictions such as regression, but many of the intuitions here also carry over to classification tasks.

- \mathcal{L} = a dataset
- $\mathbf{X}_{\mathcal{L}} = \{(y_j, \mathbf{x}_j), \quad j = 1 \dots N\}$ = the data that makes up \mathcal{L}
- $y = f(\mathbf{x}) + \epsilon$ = a noisy model from which we assume the true data is generated
- ϵ = normally distributed with mean zero and standard deviation σ_{ϵ}
- $\hat{g}_{\mathcal{L}}(\mathbf{x})$ = a predictor that we assume we have a statistical procedure (e.g. least-squares regression) for forming. The predictor gives the prediction of our model for a new data point \mathbf{x} .
- $\mathcal{C}(\mathbf{X}, \hat{g}(\mathbf{x})) = \sum_i (y_i - \hat{g}_{\mathcal{L}}(\mathbf{x}_i))^2$ = the cost function, which we have taken to be the squared error. We choose the estimator previously mentioned estimator by minimizing this.
- $\{\mathcal{L}_j\}$ = many different data sets, not just the particular training dataset \mathcal{L} that we have in hand.
- $E_{\mathcal{L}}$ = the expectation value of the cost function over $\{\mathcal{L}_j\}$ = the generalization error on all data drawn from the true model

Thus can view $\hat{g}_{\mathcal{L}}$ as a stochastic functional that depends on the dataset \mathcal{L} and can think of $E_{\mathcal{L}}$ as the expected value of the functional if we drew an infinite number of datasets $\{\mathcal{L}_1, \mathcal{L}_2, \dots\}$.

- E_{ϵ} = the expectation value over ϵ , as we would also like to average over different instances over this "noise".

Can thus decompose the expected generalization error in the following way, where in line three we use that $E[X + Y] = E[X] + E[Y]$, in line five use that $Var[X] = E[X^2] - E^2[X]$ and in line six use that our noise has zero mean ($E[\epsilon] = 0$).

$$E_{\mathcal{L}, \epsilon}[\mathcal{C}(\mathbf{X}, \hat{g}(\mathbf{x}))] = E_{\mathcal{L}, \epsilon}[\sum_i (y_i - \hat{g}_{\mathcal{L}}(\mathbf{x}_i))^2] \quad (1)$$

$$= E_{\mathcal{L}, \epsilon}[\sum_i (y_i - f(\mathbf{x}_i) + f(\mathbf{x}_i) - \hat{g}_{\mathcal{L}}(\mathbf{x}_i))^2] \quad (2)$$

$$= \sum_i E_{\epsilon}[(y_i - f(\mathbf{x}_i))^2] + E_{\mathcal{L}}[(f(\mathbf{x}_i) - \hat{g}_{\mathcal{L}}(\mathbf{x}_i))^2] + 2E_{\epsilon}[y_i - f(\mathbf{x}_i)]E_{\mathcal{L}}[f(\mathbf{x}_i) - \hat{g}_{\mathcal{L}}(\mathbf{x}_i)] \quad (3)$$

$$= \sum_i E_{\epsilon}[\epsilon^2] + E_{\mathcal{L}}[(f(\mathbf{x}_i) + 2E_{\epsilon}[\epsilon]E_{\mathcal{L}}[f(\mathbf{x}_i) - \hat{g}_{\mathcal{L}}(\mathbf{x}_i)])^2] \quad (4)$$

$$= \sum_i \sigma_{\epsilon}^2 + E_{\mathcal{L}}[(f(\mathbf{x}_i) + 2E_{\epsilon}[\epsilon]E_{\mathcal{L}}[f(\mathbf{x}_i) - \hat{g}_{\mathcal{L}}(\mathbf{x}_i)])^2] \quad (5)$$

$$= \sum_i \sigma_{\epsilon}^2 + E_{\mathcal{L}}[(f(\mathbf{x}_i) - \hat{g}_{\mathcal{L}}(\mathbf{x}_i))^2] \quad (6)$$

We further decompose the second term as

$$E_{\mathcal{L}}[(f(\mathbf{x}_i) - \hat{g}_{\mathcal{L}}(\mathbf{x}_i))^2] = E_{\mathcal{L}}[(f(\mathbf{x}_i) - E_{\mathcal{L}}[\hat{g}_{\mathcal{L}}(\mathbf{x}_i)] + E_{\mathcal{L}}[\hat{g}_{\mathcal{L}}(\mathbf{x}_i)] - \hat{g}_{\mathcal{L}}(\mathbf{x}_i))^2] \quad (7)$$

$$= E_{\mathcal{L}}[(f(\mathbf{x}_i) - E_{\mathcal{L}}[\hat{g}_{\mathcal{L}}(\mathbf{x}_i)])^2] + E_{\mathcal{L}}[(\hat{g}_{\mathcal{L}}(\mathbf{x}_i) - E_{\mathcal{L}}[\hat{g}_{\mathcal{L}}(\mathbf{x}_i)])^2] + 2E_{\mathcal{L}}[(f(\mathbf{x}_i) - E_{\mathcal{L}}[\hat{g}_{\mathcal{L}}(\mathbf{x}_i)])(\hat{g}_{\mathcal{L}}(\mathbf{x}_i) - E_{\mathcal{L}}[\hat{g}_{\mathcal{L}}(\mathbf{x}_i)])] \quad (8)$$

$$= (f(\mathbf{x}_i) - E_{\mathcal{L}}[\hat{g}_{\mathcal{L}}(\mathbf{x}_i)])^2 + E_{\mathcal{L}}[(\hat{g}_{\mathcal{L}}(\mathbf{x}_i) - E_{\mathcal{L}}[\hat{g}_{\mathcal{L}}(\mathbf{x}_i)])^2] + 2(f(\mathbf{x}_i) - E_{\mathcal{L}}[\hat{g}_{\mathcal{L}}(\mathbf{x}_i)])(E_{\mathcal{L}}[\hat{g}_{\mathcal{L}}(\mathbf{x}_i)] - E_{\mathcal{L}}[\hat{g}_{\mathcal{L}}(\mathbf{x}_i)]) \quad (9)$$

$$= (f(\mathbf{x}_i) - E_{\mathcal{L}}[\hat{g}_{\mathcal{L}}(\mathbf{x}_i)])^2 + E_{\mathcal{L}}[(\hat{g}_{\mathcal{L}}(\mathbf{x}_i) - E_{\mathcal{L}}[\hat{g}_{\mathcal{L}}(\mathbf{x}_i)])^2] + 2(f(\mathbf{x}_i) - E_{\mathcal{L}}[\hat{g}_{\mathcal{L}}(\mathbf{x}_i)])(E_{\mathcal{L}}[\hat{g}_{\mathcal{L}}(\mathbf{x}_i)] - E_{\mathcal{L}}[\hat{g}_{\mathcal{L}}(\mathbf{x}_i)]) \quad (10)$$

$$= (f(\mathbf{x}_i) - E_{\mathcal{L}}[\hat{g}_{\mathcal{L}}(\mathbf{x}_i)])^2 + E_{\mathcal{L}}[(\hat{g}_{\mathcal{L}}(\mathbf{x}_i) - E_{\mathcal{L}}[\hat{g}_{\mathcal{L}}(\mathbf{x}_i)])^2] \quad (11)$$

where we used that $(f(\mathbf{x}_i) - E_{\mathcal{L}}[\hat{g}_{\mathcal{L}}(\mathbf{x}_i)])^2$ is a constant ($f(\mathbf{x}_i)$ is deterministic, which means $E[f] = f$ and the expected value of an expected value is just that, $E[E[\hat{g}_{\mathcal{L}}]] = E[\hat{g}_{\mathcal{L}}]$) and that $E[\text{constant}] = \text{constant}$. Now we have that the two terms we are left with are called

$$Bias^2 = \sum_i (f(\mathbf{x}_i) - E_{\mathcal{L}}[\hat{g}_{\mathcal{L}}(\mathbf{x}_i)])^2 \quad (12)$$

$$Var = \sum_i E_{\mathcal{L}}[(\hat{g}_{\mathcal{L}}(\mathbf{x}_i) - E_{\mathcal{L}}[\hat{g}_{\mathcal{L}}(\mathbf{x}_i)])^2] \quad (13)$$

The bias measures the deviation of the expected value of our estimator (i. e. the asymptotic value of our estimator in the infinite data limit) from the true value. The variance measures how much our estimator fluctuates due to finite-sample effects. Combining we see that the expected out-of-sample error of our model can be decomposed as

$$E_{out} = E_{\mathcal{L}, \epsilon}[\mathcal{C}(\mathbf{X}, \hat{g}(\mathbf{x}))] = Bias^2 + Var + Noise \quad (14)$$

4 Gradient descent and its generalizations

- Ingredients shared by almost every problem in ML and data science:
 - \mathbf{X} = a dataset
 - $g(\boldsymbol{\theta})$ = a model = a function of the parameters $\boldsymbol{\theta}$
 - $\mathcal{C}(\mathbf{X}, g(\boldsymbol{\theta}))$ = a cost function allowing us to judge how well the model explains the observations. We fit the model by finding the values of $\boldsymbol{\theta}$ that minimize the cost function.

Here we discuss one of the most powerful and widely used classes of methods for performing this minimization - gradient descent and its generalizations. Main idea: iteratively adjust the parameters in the direction where the gradient of the cost function is large and negative. In this way the training procedure ensures the parameters flow towards a *local* minimum of the cost function.

- Underlying reason training a LM algo is hard is the cost functions we wish to optimize are usually complicated, rugged, non-convex functions in a high-dimensional space with many local minima.
- Gradient descent (GD) and Newton's method:
 - $E(\boldsymbol{\theta})$ = function we wish to minimize. "Energy function". In ML context $E(\boldsymbol{\theta}) = \mathcal{C}(\mathbf{X}, g(\boldsymbol{\theta}))$.
 - Almost always: $E(\boldsymbol{\theta}) = \sum_{i=1}^n e_i(\mathbf{x}_i, \boldsymbol{\theta})$ = a sum over n data points. Fex: Linear regression: e_i = the mean square error for data point i . Logistic regression: e_i = the cross-entropy. Call e_i the energy function to make analogy to physical systems.
- Simplest GD: start with an initial value $\boldsymbol{\theta}_0$, then update according to

$$\mathbf{v}_t = \eta_t \nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta}_t) \quad (15)$$

$$\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t - \mathbf{v}_t \quad (16)$$

where we have introduced the *learning rate*, η_t , that controls how big a step we should take in the direction of the gradient at time t . For sufficiently small η_t this method will converge to a *local minimum* of the cost func. But a small η_t comes at a computational cost. If it's smaller we need more steps to reach the minimum. But if it's too large we can overshoot the minimum and the algo becomes unstable (either oscillates or even moves away from the minimum). See fig 7. In practice, one usually specifies a "schedule" that decreases η_t at long times. Common schedules include power law and exponential decays in time.

- Contrast with Newton’s method to better understand this behavior of GD. In Newton’s method, we choose the step \mathbf{v} for the parameters in such a way as to minimize a second-order Taylor expansion to the energy function

$$E(\boldsymbol{\theta} + \mathbf{v}) \approx E(\boldsymbol{\theta}) + \nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta}) \mathbf{v} + \frac{1}{2} \mathbf{v}^T H(\boldsymbol{\theta}) \mathbf{v} \quad (17)$$

where $H(\boldsymbol{\theta})$ = the Hessian matrix of second derivatives. Differentiating this equation wrt \mathbf{v} and noting that for the optimal value \mathbf{v}_{opt} we expect $\nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta} + \mathbf{v}_{opt}) = 0$, yields the equation

$$0 = \nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta}) + H(\boldsymbol{\theta}) \mathbf{v}_{opt} \quad (18)$$

Rearranging results in the desired update rules for Newton’s method

$$\mathbf{v}_t = H^{-1}(\boldsymbol{\theta}_t) \nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta}_t) \quad (19)$$

$$\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t - \mathbf{v}_t \quad (20)$$

Have no guarantee the Hessian is well conditioned, so often replaces the Hessian inverse $H^{-1}(\boldsymbol{\theta}_t)$ by some suitably regularized pseudo-inverse such as $[H(\boldsymbol{\theta}_t) + \epsilon I]^{-1}$ w/a small ϵ parameter.

- For ML, Newton’s method not practical for two interrelated reasons:
 - Calculating a Hessian is an extremely expensive numerical computation.
 - Even if we employ first order approximation methods to approximate the Hessian (commonly called quasi-Newton methods), we must store and invert a matrix with n^2 entries, n =the number of parameters. For models with millions of parameters such as those commonly employed in neural networks, this is close to impossible with present-day computational power.
- Important intuition from Newton’s method to modify GD: Netwon’s method automatically adapts the learning rate of different parameters depending on the Hessian matrix. Whereas simple GD has the same learning rate for all the parameters. The Hessian encodes the curvature of the surface we’re minimizing. Specifically, the singular values of the Hessian are inversly proportional to the squares of the local curvatures of the surface.

$$\text{the Hessian's singular values} \propto \frac{1}{(\text{the local curvatures of the surface})^2} \quad (21)$$

Newton’s method thus automatically adjusts the step size so that one takes larger steps in flat directions with small curvatures and smaller steps in steep directions with large curvature.

- Consider special case: Using GD to find minimum of a quadratic energy func of a single parameter θ . Given current value of θ , find η_{opt} = the η that lets us reach the minimum in a single step. To find it expand the energy func to second order around the current value

$$E(\theta + v) = E(\theta_c) + \partial_\theta E(\theta)v + \frac{1}{2}\partial_\theta^2 E(\theta)v^2 \quad (22)$$

We want to find the step v such that $\theta + v$ is a stationary point. That is we seek to solve the equation that sets the derivative of this last expression wrt v equal to zero:

$$0 = \partial_\theta E(\theta) + \partial_\theta^2 E(\theta)v \quad (23)$$

$$v = -\partial_\theta E(\theta)[\partial_\theta^2 E(\theta)]^{-1} \quad (24)$$

$$\Rightarrow \theta_{min} = \theta - v \quad (25)$$

$$\theta_{min} = \theta - [\partial_\theta^2 E(\theta)]^{-1}\partial_\theta E(\theta) \quad (26)$$

Comparing with the previously outlined GD update rule tells us that

$$\eta_{opt} = [\partial_\theta^2 E(\theta)]^{-1} \quad (27)$$

Four qualitatively different regimes possible (fig 8):

- $\eta < \eta_{opt}$: GD will take multiple small steps to reach the bottom of the potential.
- $\eta = \eta_{opt}$: GD reaches the bottom of the potential in a single step
- $\eta_{opt} < \eta < 2\eta_{opt}$: GD oscillates across both sides of the potential before eventually converging to the minima.
- $\eta > 2\eta_{opt}$: GD diverges!
- Straightforward to generalize to the multidimensional case: The natural multidimensional generalization of the second derivative = the Hessian $H(\theta)$. Can always perform a singular value decomposition (= a rotation bt an orthogonal matrix for quadratic minima where the Hessian is symmetric) and consider the Hessian's singular values $\{\lambda\}$. If we use a single learning rate for all parameters, in analogy with the η_{opt} found above, convergence requires that

$$\eta < \frac{2}{\lambda_{max}} \quad (28)$$

where λ_{max} =the Hessian's largest singular value. If the minimum λ_{min} differs significantly from λ_{max} , then convergence in the λ_{min} direction will be extremely slow. Can show that convergence time scales w/ the condition number

$$\kappa = \frac{\lambda_{max}}{\lambda_{min}} \quad (29)$$

- Limitations of the simplest GD algo:
 - *GD finds local minima* Since GD is deterministic it converges to a local minimum of our energy func. May lead to poor performance in ML. A similar problem is encountered in physics and overcome by methods like **simulated annealing** that introduce a fictitious "temperature" which is eventually taken to zero. The "temperature" term introduces stochasticity in the form of thermal fluctuations that allow the algo to thermally tunnel over energy barriers. Suggests we should in ML modify GD to include stochasticity.
 - *GD is sensitive to initial conditions* Initial conditions matter, as a consequence of GD's local nature. Thus, very important to think about how you initialize, both for simple and more complicated GD variant introduced later.
 - *Gradients are computationally expensive to calculate for large datasets* As previously mentioned, in many statistics and ML cases, the energy func is a sum of terms, one for each data point = to calc the gradient we must sum over all n data points. Do this at each GD point = extremely computationally expensive. An ingenious solution: calc the gradients using small subsets of data = "mini batches". This also introduces stochasticity into our algo.
 - *GD very sensitive to choices of learning rates* Small lr = slow, large lr = possible divergence/poor results. Also, depending on the local landscape, we have to modify lr to ensure convergence. Ideally would "adaptively" choose lr to match the landscape.
 - *GD treats all directions in parameter space uniformly* Unlike Newton's method. Thus the maximum lr is set by the behavior of the steepest direction and this can significantly slow training. Would ideally like large steps in flat dir and small steps in steep dir. Since we are exploring rugged landscapes where curvatures change, this requires us to keep track of not only the gradient but second derivatives of the energy func (as discussed calc Hessian would be ideal, but proves too computationally expensive).
 - *GD can take exponential time to escape saddle points, even with random initialization* As mentioned extremely sensitive to initial conditions since it determines the particular local minimum GD will reach. But, even with a good initialization scheme (through the introduction of randomness) GD can still take exponential time to escape saddle points, prevalent in high-dimensional space, even for non-pathological objective functions. **There are modified GD methods developed recently to accelerate the escape, see reference.**

- Stochastic gradient descent (SGD) with mini-batches: Stochasticity is incorporated by approximating the gradient on a subset of the data called a mini-batch.
 - Size of mini-batch almost always \ll total number of data points n .
 - Typical mini-batch sizes ranging from ten to a few hundred data points
 - If there are n points, and the mini-batch size is M , there will be n/M mini-batches.
 - Denote these mini-batches by B_k where $k = 1 \dots n/M$.

Thus in SGD at each GD step we approximate the gradient using a minibatch B_k ,

$$\nabla_{\theta} = \sum_i^n \nabla_{\theta} e_i(\mathbf{x}_i, \theta) \rightarrow \sum_{i \in B_k} \nabla_{\theta} e_i(\mathbf{x}_i, \theta) \quad (30)$$

cycling over all M minibatches. A full iteration over all n data points = using all M minibatches = an *epoch*. Denote the minibatch approximation to the gradient by

$$\nabla_{\theta} E^{MB}(\theta) = \sum_{i \in B_k}^M \nabla_{\theta} e_i(\mathbf{x}_i, \theta) \quad (31)$$

Then SGD algo is

$$\mathbf{v}_t = \eta_t \nabla_{\theta} E^{MB}(\theta) \quad (32)$$

$$\theta_{t+1} = \theta_t - \mathbf{v}_t \quad (33)$$

Two important benefits to SGD.

- Introduces stochasticity and decreases chance the fitting algo gets stuck in isolated local minima.
- Significantly speeds up the calc as one does not have to use all n data points to calc the gradient.
- Empirical and theoretical work suggests SGD has additional benefits - one significant being that introducing stochasticity is thought to act as a **natural regularizer that prevents overfitting** in deep, isolated minima.
- Adding momentum: SGD almost always used with a "momentum"/inertia term serving as a memory of the direction we are moving in parameter space. Implemented as

$$\mathbf{v}_t = \gamma \mathbf{v}_{t-1} + \eta_t \nabla_{\theta} E(\theta_t) \quad (34)$$

$$\theta_{t+1} = \theta_t - \mathbf{v}_t \quad (35)$$

where γ =a momentum parameter with $0 \leq \gamma \leq 1$ (and have dropped the explicit notation indicating the gradient is taken over a minibatch). This is gradient descent with momentum (GDM). Clear that

- \mathbf{v}_t is a running average
- $(1 - \gamma)^{-1}$ sets the characteristic time scale for the memory used in the averaging procedure.
- $\gamma = 0$ reduces down to ordinary SGD

Equivalent way of writing updates:

$$\Delta\theta_{t+1} = \gamma\Delta\theta_t - \eta_t \nabla_{\theta} E(\theta_t) \quad (36)$$

where $\Delta\theta_t = \theta_t - \theta_{t-1}$.

Getting intuition: Consider simple physical analogy with a particle of mass m moving in a viscous medium with drag coefficient μ and potential $E(\mathbf{w})$, \mathbf{w} =the particle's position. It's motion then described by

$$m \frac{d^2 \mathbf{w}}{dt^2} + \mu \frac{d\mathbf{w}}{dt} = -\nabla_{\mathbf{w}} E(\mathbf{w}) \quad (37)$$

Discretize in usual way to get

$$m \frac{\mathbf{w}_{t+\Delta t} - 2\mathbf{w}_t + \mathbf{w}_{t-\Delta t}}{(\Delta t)^2} + \mu \frac{\mathbf{w}_{t+\Delta t} - \mathbf{w}_t}{\Delta t} = -\nabla_{\mathbf{w}} E(\mathbf{w}) \quad (38)$$

Rearrange to write as

$$\Delta\mathbf{w}_{t+\Delta t} = -\frac{(\Delta t)^2}{m + \mu\Delta t} \nabla_{\mathbf{w}} E(\mathbf{w}) + \frac{m}{m + \mu\Delta t} \Delta\mathbf{w}_t \quad (39)$$

Notice it's identical to our GDM update rule defined above. We may thus identify the momentum parameter and learning rate with the mass of the particle and the viscous drag:

$$\gamma = \frac{m}{m + \mu\Delta t} \quad (40)$$

$$\eta = \frac{(\Delta t)^2}{m + \mu\Delta t} \quad (41)$$

Thus as suggested by the name the momentum parameter is proportional to the mass of the particle and effectively provides inertia. Also, in the large viscosity/small learning rate limit, our memory scales as $(1 - \gamma)^{-1} \approx m/(\mu\Delta t)$

Why momentum useful?

- Helps the algo gain speed in directions with persistent but small gradients even in the presence of stochasticity, while suppressing oscillations in high-curvature directions.
- Has been argued first-order methods (with appropriate initial conditions) can perform comparable to more expensive second-order methods, especially in context of complex deep learning models (reference)
- Studies suggest benefits of momentum especially pronounced in complex models in the initial "transient phase" of training, rather than during subsequent fine-tuning of a coarse minimum. Because in the transient phase, correlations in the gradient persist across many GD steps, accentuating the role of inertia and memory.

These can be even more pronounced using a slightly modified algo, Nesterov Accelerated Gradient (NAG). Rather than calc gradient at the current parameters $\nabla_{\theta}E(\theta_t)$, calc the gradient at the expected value of the parameters given our current momentum $\nabla_{\theta}E(\theta_t + \gamma \mathbf{v}_{t-1})$. One major advantage is it allows for larger learning rate than GDM for same choice of γ .

- Methods that use the second moment of the gradient: In SGD, with and without momentum, we still have to specify a "schedule" for tuning the learning rates η_t as a func of time. As touched upon before this presents dilemmas.
 - The lr is limited by the steepest direction which can change depending on the current position in the landscape. To circumvent this, ideally our algo would keep track of curvature and take large steps in shallow/flat dirs and small steps in steep/narrow dirs.
 - Second-order methods accomplish this by calc or approximating the Hessian and normalizing the lr by the curvature.
 - But this is very computationally expensive for extremely large models.
 - Ideally, we would like to adaptively change the step size to match the landscape without paying the steep computational price of calculating/approximating the Hessian.
 - Recently introduced a number of methods that accomplish this by tracking not only the gradient, but also the second moment of the gradient. These include AdaGrad, AdaDelta, RMS-Prop, and ADAM. We'll discuss the last two.

RMS prop: In addition to keeping a running average of the first moment of the gradient, also keep track of the second moment denoted

by $\mathbf{s}_t = \mathbb{E}[\mathbf{g}_t^2]$. Update rule is

$$\mathbf{g}_t = \nabla_{\theta} E(\boldsymbol{\theta}) \quad (42)$$

$$\mathbf{s}_t = \beta \mathbf{s}_{t-1} + (1 - \beta) \mathbf{g}_t^2 \quad (43)$$

$$\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t - \eta_t \frac{\mathbf{g}_t}{\sqrt{\mathbf{s}_t + \epsilon}} \quad (44)$$

where

- β = controls the averaging time of the second moment and is typically $\beta = 0.9$
- η_t = a learning rate, typically 10^{-3}
- $\epsilon \sim 10^{-8}$ = a small regularization constant to prevent divergences
- Multiplication and division by vectors is understood as element-wise operations
- Clear from the formula that lr reduced in dirs where the norm of the gradient is consistently large. This greatly speeds up the convergence by allowing us to use a larger lr for flat dirs.

The ADAM optimizer: Keep a running average of both the first and second moment of the gradient ($\mathbf{m}_t = \mathbb{E}[\mathbf{g}_t]$ and $\mathbf{s}_t = \mathbb{E}[\mathbf{g}_t^2]$ respectively) - use this info to adaptively change the lr for different parameters. Also performs an additional bias correction to account for the fact that we're estimating the first two moments of the gradient using a running average (denoted below by the hats). Update rule is (multiplication and division by vectors again understood to be element-wise):

$$\mathbf{g}_t = \nabla_{\theta} E(\boldsymbol{\theta}) \quad (45)$$

$$\mathbf{m}_t = \beta_1 \mathbf{m}_{t-1} + (1 - \beta_1) \mathbf{g}_t \quad (46)$$

$$\mathbf{s}_t = \beta_2 \mathbf{s}_{t-1} + (1 - \beta_2) \mathbf{g}_t^2 \quad (47)$$

$$\hat{\mathbf{m}}_t = \frac{\mathbf{m}_t}{1 - \beta_1^t} \quad (48)$$

$$\hat{\mathbf{s}}_t = \frac{\mathbf{s}_t}{1 - \beta_2^t} \quad (49)$$

$$\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t - \eta_t \frac{\hat{\mathbf{m}}_t}{\sqrt{\hat{\mathbf{s}}_t + \epsilon}} \quad (50)$$

where

- β_1 and β_2 set the memory lifetime of the first and second moment, typically taken as 0.9 and 0.99 respectively
- η and ϵ identical to RMS prop.

Like in RMSprop the effective step size of a parameter depends on the magnitude of its gradient squared. To better understand, let's rewrite this expression in terms of the variance $\sigma_t^2 = \hat{\mathbf{s}}_t - (\hat{\mathbf{m}}_t)^2$. Consider a single parameter θ_t . Update rule of it is given by

$$\Delta\theta_{t+1} = -\eta_t \frac{\hat{\mathbf{m}}_t}{\sqrt{\sigma_t^2 + \hat{\mathbf{s}}_t^2 + \epsilon}} \quad (51)$$

We examine the limiting cases of this expression.

- Assume our gradient estimates are consistent so the variance is small. Then update rule tends to $\Delta\theta_{t+1} \rightarrow -\eta_t$ (assumed $\hat{\mathbf{m}}_t \gg \epsilon$). This is equivalent to cutting off large persistent gradients at 1 and limiting the max step size in steep directions.
- Imagine the gradient is widely fluctuating between GD steps. Then $\sigma^2 \gg \hat{\mathbf{m}}_t^2$ so our update becomes $\Delta\theta_{t+1} \rightarrow -\eta_t \hat{\mathbf{m}}_t / \sigma_t$. AKA, we adapt our lr so that

learning rate \propto signal-to-noise ratio (i.e. the mean in units of the standard deviation)

- = extremely desirable from physical point of view. The std serves as a natural adaptive scale for deciding whether a gradient is large or small.
- Thus, ADAM has the beneficial effects of adapting our step size so that we cut off large gradient dirs (and hence prevent oscillations and divergences) and measuring gradients in terms of a natural length scale, the std σ_t .
- Above discussion also explains empirical observations showing that the performance of both ADAM and RMSprop is drastically reduced if the square root omitted in the update rule.
- Also worth noting recent studies have shown adaptive methods like RMSprop, ADAM, and AdaGrad tend to generalize worse than SGD in classification tasks, though they achieve smaller training error. See refs.
- Comparisons of various methods: Visualize the performance of the 5 discussed methods GD, GDM, NAG, ADAM and RMSprop by using Beale's function:

$$f(x, y) = (1.5 - x + xy)^2 + (2.25 - x + xy^2)^2 + (2.625 - x + xy^3)^2 \quad (52)$$

It has a global minimum at $(x, y) = (3, 0.5)$ and an interesting structure. Fig 9:

- Shows the results of using all five methods for $N_{steps} = 10^4$ steps from three different initial conditions.
- GD, GDM, NAG learning rate: $\eta = 10^{-6}$

- RMSprop, ADAM learning rate: $\eta = 10^{-3}$ (can be higher due to their adaptive step sizes. thus these methods tend to be quicker at navigating the landscape.)
- Notice in some cases (e.g. initial condition $(-1, 4)$), the trajectories do not find the global min but instead follow the deep, narrow ravine that occurs along $y = 1$. This kind of landscape structure is generic in high-dimensional spaces where saddle points proliferate.
- Gradient descent in practice: practical tips (especially in the context of deep neural networks):
 - *Randomize the data when making mini-batches* Important to randomly shuffle the data when forming the batches. Otherwise GD can fit spurious correlations resulting from the order in which data is presented.
 - *Transform your inputs* As discussed learning=difficult when landscape has mix of steep and flat dirs. Trick: standardize the data by subtracting the mean and normalizing the variance of input variables. Whenever possible also decorrelate the inputs. Why helpful - consider the case of linear regression: for the squared error cost func, the Hessian of the energy matrix = the correlation matrix between inputs. = by standardizing the inputs we ensure the landscape looks homogeneous in all dirs in parameter space. Most deep networks can be viewed as linear transformations followed by a non-linearity at each layer, thus we expect this intuition to hold beyond the linear case.
 - *Monitor the out-of-sample performance* Always monitor the performance on a validation set = a small portion of the training data kept out of the training process to serve as a proxy for the test set. Validation error starting to increase = model being overfit. Terminate the learning. This *early stopping* significantly improves performance in many settings.
 - *Adaptive optimization methods don't always have good generalization* As mentioned recent studies have shown adaptive methods such as ADAM, RMSprop and AdaGrad to have poor generalization compared to SGD or GDM, particularly in the high-dimensional limit = the number of parameters exceeds the number of data points. Although not clear why these methods perform so well in training deep neural networks such as generative adversarial networks (GANs), simpler procedures like properly tuned SGD may work as well or better in these applications.

5 Overview of Bayesian inference

Statistical modeling focus: estimation/prediction of unknown quantities. Bayesian methods premise: probability can be used as mathematical tool for describing uncertainty. Similar in spirit to physics statistical mechanics - where we use probability to describe the behavior of large systems where we cannot know the positions and momenta of all particles even if the system itself is fully deterministic (at least classically). This section gives introduction to Bayesian inference, w/special emphasis on its logic=Bayesian reasoning and connections to ML.

5.1 Bayes rule

Must specify two functions

- $p(\mathbf{X}|\mathbf{w})$ = the likelihood function, which describes the probability of observing a dataset \mathbf{X} for a given value of the unknown parameters \mathbf{w} . The func should be considered a func of the parameters \mathbf{w} with the data \mathbf{X} held fixed.
- $p(\mathbf{w})$ = the prior distribution, which describes any knowledge we have about the parameters before we collect the data.
- with these two we can compute the posterior distribution via Baye's rule:

$$p(\mathbf{w}|\mathbf{X}) = \frac{p(\mathbf{X}|\mathbf{w})p(\mathbf{w})}{\int d\mathbf{w} p(\mathbf{X}|\mathbf{w})p(\mathbf{w})} \quad (53)$$

It describes our knowledge about the unknown parameter \mathbf{w} after observing the data \mathbf{X} . In many casing computing the normalizing constant (i.e. the partition function $p(\mathbf{X}) = \int d\mathbf{w} p(\mathbf{X}|\mathbf{w})p(\mathbf{w})$) and Markov Chain Monte Carlo (MCMC) methods are needed to draw random samples from the posterior dist.

The likelihood func = common feature in both classical statistics and Bayesian inference. Determined by the model and the measurement noise. Many statistical procedures fex least-square fitting can be cast into the formalism of *Maximum Likelihood Estimation* (MLE). In MLE one chooses the parameters $\hat{\mathbf{w}}$ that maximize the likelihood (or equivalently the log-likelihood since log a monotonic func) of the observed data:

$$\hat{\mathbf{w}} = \operatorname{argmax}_w \{ \log(p(\mathbf{X}|\mathbf{w})) \} \quad (54)$$

Aka, in MLE we choose the parameters that maximize the probability of seeing the observed data given our generative model. MLE is an important concept in both frequentist and Bayesian statistics.

The prior distribution, by contrast, is uniquely Bayesian. Two general classes of priors:

- If we do not have any specialized knowledge of \mathbf{w} before we look at the data we would like to select an *uninformative* prior that reflects our ignorance
- Otherwise we select an *informative* prior that accurately reflects the knowledge we have about \mathbf{w} .

There is a large literature on uninformative priors, including reparametrization invariant priors that would be of interest to physicists, but here we'll focus on *informative* priors. Using this tends to decrease the variance of the posterior distribution while, potentially, increasing its bias. Beneficial if decrease in variance \gg increase in bias.

In high-dimensional problems it's reasonable to assume many of the parameters not strongly relevant = many of the parameters will be zero or close to zero. Can express this belief using two commonly used priors:

- $p(\mathbf{w}|\lambda) = \prod_j \sqrt{\frac{\lambda}{2\pi}} e^{-\lambda w_j^2}$ = the Gaussian prior = used to express the assumption many of the parameters will be small
- $p(\mathbf{w}|\lambda) = \prod_j \frac{\lambda}{2} e^{-\lambda |w_j|}$ = the Laplace prior = used to express the assumption that many of the parameters will be zero.

Come back to this in section VI.F.

5.2 Bayesian decisions

We have seen how to compute the posterior distribution, which expresses our knowledge about the parameters \mathbf{w} . Mostly however we need to summarize our knowledge and pick a single "best" value for them. In principle the specific value should be chosen to maximize a utility function. In practice however, usually use one of two choices:

- $\langle \mathbf{w} \rangle = \int d\mathbf{w} \mathbf{w} p(\mathbf{w}|\mathbf{X})$ = the posterior mean = the Bayes estimate (minimizes the mean-squared error)
- $\hat{\mathbf{w}} = \operatorname{argmax}_{\mathbf{w}} p(\mathbf{w}|\mathbf{X})$ = the posterior mode = the maximum-a-posteriori = MAP estimate (easier to compute)

5.3 Hyperparameters

The Gaussian and Laplace prior dist. both have an extra parameter λ = *hyperparameter* = *nuisance variable*. Has to be chosen somehow.

- One standard approach: define another prior dist for λ , usually using an uninformative prior, and average the posterior dist over all choices of λ . = a hierarchical prior. However computing averages often requires long MCMC simulations that are computationally expensive.

- Simpler: find a good λ value using an optimization procedure. Will discuss later.

6 Linear regression

In this section we take a closer look at the ideas of the optimal choice of predictor depending on the choice of fitted function and underlying noise level, and model complexity, the bias-variance decomposition, the statistical meaning of learning. As previously fitting a given set of samples (y_i, \mathbf{x}_i) means relating the independent variables \mathbf{x}_i to their responses y_i .

Formulating the problem:

- $\mathcal{D} = \{(y_i, \mathbf{x}_i)\}_{i=1}^n$ = a given dataset with n samples, where \mathbf{x}_i = the i -th observation vector while y_i = its corresponding (scalar) response.
- $\mathbf{x}_i \in \mathbb{R}^p$ = assume every sample has p features.
- f = the true function/model that generated these samples via $y_i = f(\mathbf{x}_i; \mathbf{w}_{true}) + \epsilon_i = \mathbf{x}_i^T \mathbf{w}_{true} + \epsilon_i$ for some unknown but fixed $\mathbf{w}_{true} \in \mathbb{R}^p$.
- g = the function we wish to find, with parameters \mathbf{w} fit to the data \mathcal{D} , that can best approximate f . When we have a $\hat{\mathbf{w}}$ such that $g(\mathbf{x}; \hat{\mathbf{w}})$ yields our best estimate of f , we can use this g to make predictions about the response y_0 for a new data point \mathbf{x}_0 .
- L^p for any real number $p \geq 1$ = the norm of a vector $\mathbf{x} = (x_1, \dots, x_d) \in \mathbb{R}^d$, defined as

$$\|\mathbf{x}\|_p = (|x_1|^p + \dots + |x_d|^p)^{\frac{1}{p}} \quad (55)$$

6.1 Least-square regression

Ordinary least squares linear regression (OLS) = the minimization of the L_2 norm of the difference between the response y_i and the predictor $g(\mathbf{x}_i; \mathbf{w}) = \mathbf{x}_i^T \mathbf{w}$:

$$\min_{\mathbf{w} \in \mathbb{R}^p} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2 = \min_{\mathbf{w} \in \mathbb{R}^p} \sum_{i=1}^n (\mathbf{x}_i^T \mathbf{w} - y_i)^2 \quad (56)$$

Geometrically, the predictor func g defines a hyperplane in \mathbb{R}^p . Thus minimizing least-square error = minimizing the sum of all projections (=residuals) for all points \mathbf{x}_i to this hyperplane (fig 10). Denote the solution to this as $\hat{\mathbf{w}}_{LS}$:

$$\hat{\mathbf{w}}_{LS} = \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^p} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2 \quad (57)$$

which after differentiation leads to

$$\hat{\mathbf{w}}_{LS} = (X^T X)^{-1} X^T \mathbf{y} \quad (58)$$

where we assumed $X^T X$ is invertible.

- Often the case when $n \gg p$. Formally: if $\text{rank}(X) = p$, namely, the predictors X_1, \dots, X_p (=the columns of X) are linearly independent, then $\hat{\mathbf{w}}_{LS}$ is unique.
- $\text{rank}(X) < p$: happens when $p > n$. $X^T X$ singular. Implying there are infinitely many solutions to the least squares problem. Can then show that if \mathbf{w}_0 a solution, then also $\mathbf{w}_0 + \eta$ also a solution for any η which satisfies $\mathbf{X}\eta = 0$ ($= \eta \in \text{null}(X)$).

Having the solution, we can calculate \mathbf{y} , the best fit of our data X , as

$$\hat{\mathbf{y}} = X \hat{\mathbf{w}}_{LS} = P_X \mathbf{y} \quad (59)$$

$$\text{where } P_X = X(X^T X)^{-1} X^T \quad (60)$$

Geometrically, P_X = the projection matrix which acts on \mathbf{y} and projects it onto the column space of X , which is spanned by the predictors X_1, \dots, X_p (fig 11).

Notice: we found the optimal solution $\hat{\mathbf{w}}_{LS}$ in one shot, no iterative optimization.

Have explained the difference between *learning* and *fitting* lies in the prediction on "unseen" data. Must thus examine out-of-sample error. Following our previos definitions in section 3 of \bar{E}_{in} and \bar{E}_{out} , the average in-sample and out-of-sample error can be shown to be

$$\bar{E}_{in} = \sigma^2 \left(1 - \frac{p}{n}\right) \quad (61)$$

$$\bar{E}_{out} = \sigma^2 \left(1 + \frac{p}{n}\right) \quad (62)$$

provided we obtain $\hat{\mathbf{w}}_{LS}$ from i.i.d. samples X and \mathbf{y} generated through $\mathbf{y} = X \mathbf{w}_{true} + \epsilon$ (this requires ϵ is a noise vector whose elements are i.i.d. of 0 mean nad variance σ^2 , and is independent of the samples X). Can thus calc the average generalization error explicitly:

$$|\bar{E}_{in} - \bar{E}_{out}| = 2\sigma^2 \frac{p}{n} \quad (63)$$

This imparts important message:

- $p \gg n$ (= high dimensional data): the generalization error is extremely large = the model is not learning
- $p \approx n$: even now we might still not learn well due to the intrinsic noise σ^2 . One way of amelioration = reularization. Mainly two forms:
 - *Ridge regression*: employs an L_2 penalty
 - *LASSO*: uses an L_1 penalty

6.2 Ridge-regression

We here study effect of adding to the least squares loss function a *regularizer* defined as the L_2 norm of the parameter vector we want to optimize over = wish to solve the following *penalized* regression problem, *Ridge regression*:

$$\hat{\mathbf{w}}_{Ridge}(\lambda) = \underset{\mathbf{w} \in \mathbb{R}^p}{\operatorname{argmin}} (||X\mathbf{w} - \mathbf{y}||_2^2 + \lambda ||\mathbf{w}||_2^2) \quad (64)$$

Equivalent to the following *constrained* optimization problem

$$\hat{\mathbf{w}}_{Ridge}(t) = \underset{\mathbf{w} \in \mathbb{R}^p: ||\mathbf{w}||_2^2 \leq t}{\operatorname{argmin}} ||X\mathbf{w} - \mathbf{y}||_2^2 \quad (65)$$

Means for any $t \geq 0$ and $\hat{\mathbf{w}}_{Ridge}$ in the last equation there exists a $\lambda \geq 0$ such that $\hat{\mathbf{w}}_{Ridge}$ solves the second last equation, and vice versa (note: equivalence between the penalized and the constrained (regularized) form of least square does not always hold. It holds for Ridge and LASSO, but not for best subset selection defined by choosing a L^0 norm).

With this equivalence it's obvious that by adding a regularization term $\lambda ||\mathbf{w}||_2^2$ to our ls loss func, we are effectively constraining the magnitude of the parameter vector learned from the data. To see this, let's evaluate the second last equation explicitly. Differentiating w.r.t \mathbf{w} gives

$$\hat{\mathbf{w}}_{Ridge}(\lambda) = (X^T X + \lambda I_{p \times p})^{-1} X^T \mathbf{y} \quad (66)$$

$$\text{If } X \text{ orthogonal:} \quad = \frac{\hat{\mathbf{w}}_{LS}}{1 + \lambda} \quad (67)$$

Implies the ridge estimate is merely the least squares estimate scaled by a factor $(1 + \lambda)^{-1}$.

Can we derive similar relation between the fitted vector $\hat{\mathbf{y}} = X\hat{\mathbf{w}}_{Ridge}$ and the prediction made by ls linear regression? We do a singular value decomposition (SVD) on X . The SVD of an $n \times p$ matrix X has the form

$$X = UDV^T \quad (68)$$

where $[U]_{n \times p}$ and $[V]_{p \times p}$ are orthogonal matrices such that the columns of U span the column space of X while the columns of V span the row space of X . $[D]_{p \times p} = \operatorname{diag}(d_1, d_2, \dots, d_p)$ is a diagonal matrix with entries $d_1 \geq d_2 \geq \dots d_p \geq 0$ called the singular values of X . X is singular if there is at least one $d_j = 0$. By writing X in terms of its SVD, one can recast the Ridge estimator as

$$\hat{\mathbf{w}}_{Ridge} = V(D^2 + \lambda I)^{-1} D U^T \mathbf{y} \quad (69)$$

which implies that the Ridge predictor satisfies

$$\hat{\mathbf{y}}_{Ridge} = X \hat{\mathbf{w}}_{Ridge} \quad (70)$$

$$= U D (D^2 + \lambda I)^{-1} D U^T \mathbf{y} \quad (71)$$

$$= \sum_{j=1}^p \mathbf{u}_j \frac{d_j^2}{d_j^2 + \lambda} \mathbf{u}_j^T \mathbf{y} \quad (72)$$

$$\leq U U^T \mathbf{y} \quad (73)$$

$$= X \hat{\mathbf{y}} \equiv \hat{\mathbf{y}}_{LS} \quad (74)$$

where in the inequality step we used SVD to simplify $\hat{\mathbf{w}}_{LS}$ and asumed $\lambda \geq 0$. Comparing the third last line and last line it's clear that

- In order to compute the fitted vector $\hat{\mathbf{y}}$, both Ridge and ls linear regression have to project \mathbf{y} to the column space of X .
- Only difference: Ridge regression further shrinks each basis component j by a factor $d_j^2/(d_j^2 + \lambda)$

6.3 LASSO and sparse regression

Here study the effects of adding an L_1 regularization penalty, called *LASSO* (=least absolute shrinkage and selection operator). LASSO in the penalized form is defined by the following regularized regression problem:

$$\hat{\mathbf{w}}_{LASSO}(\lambda) = \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^p} \|X\mathbf{w} - \mathbf{y}\|_2^2 + \lambda \|\mathbf{w}\|_1 \quad (75)$$

As with Ridge there is another formulation based on constrained optimization:

$$\hat{\mathbf{w}}_{LASSO}(t) = \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^p: \|\mathbf{w}\|_1 \leq t} \|X\mathbf{w} - \mathbf{y}\|_2^2 \quad (76)$$

- The equivalence interpretation same as in Ridge, namely: for any $t \geq 0$ and solution $\hat{\mathbf{w}}_{LASSO}$ in the last equation there is a $\lambda \geq 0$ such that $\hat{\mathbf{w}}_{LASSO}$ solves the second last equation. And vice versa.
- However, to get the analytical solution of LASSO cannot simply take gradient wrt \mathbf{w} , since the L_1 -regularizer is not everywhere differentiable, in particular at any point where $w_j = 0$ (fig 13).
- Nonetheless, LASSO is a **convex** problem. Can thus invoke the *sub-gradient optimality condition* in optimization theory to obtain solution.

For simple notation, only show the solution for when X is orthogonal:

$$X \text{ orthogonal: } \hat{w}_j^{LASSO}(\lambda) = \operatorname{sign}(\hat{w}_j^{LS})(|\hat{w}_j^{LS}| - \lambda)_+ \quad (77)$$

where $(x)_+ =$ the positive part of x .

Fig 12 compares the Ridge and LASSO solutions.

- As mentioned, Ridge is the LS solution scaled by a factor $(1 + \lambda)$. LASSO does something called "soft-thresholding".
- LASSO tends to give sparse solutions, meaning many components of $\hat{\mathbf{w}}_{LASSO}$ are zero. Fig 13 gives intuitive explanation of this. In short:
 - To solve a constrained optimization problem with a fixed regularization strength $t \geq 0$ one first carves out the "feasible region" specified by the regularizer in the $\{w_1, \dots, w_d\}$ space = a solution $\hat{\mathbf{w}}_0$ only legitimate if falls within this region.
 - Then one plots the contours of the ls regressors in an increasing manner until the contour touches the feasible region.
 - The point where this occurs is the solution to our optimization problem.
 - Loosely speaking, since the L_1 regularizer of LASSO has sharp protrusions(=vertices) along the axes, and because the regressor contours are in the shape of ovals (it's quadratic in \mathbf{w}), their intersection tends to occur at the vertex of the feasibility region = implying the solution vector will be sparse. (the vertices correspond to parameter vectors \mathbf{w} with only one non-vanishing component).

See fig 14 and 15 for comparison of Ridge and Lasso on Diabetes dataset.

6.4 Using linear regression to learn the Ising Hamiltonian

Task: learning the Hamiltonian for the Ising model. Given an ensemble of random spin configurations, and assigned to each state its energy, generated from the 1D Ising model:

$$H = -J \sum_{j=1}^L S_j S_{j+1} \quad (78)$$

where

- J = the nearest-neighbour spin interaction
- $S_j \in \{\pm 1\}$ = a spin variable
- Assume data was generated with $J = 1$
- You are handed the dataset $\mathcal{D} = (\{S_j\}_{j=1}^L, E_j)$ without knowing what the E_j number mean
- The $\{S_j\}_{j=1}^L$ configuration can be interpreted in many ways: outcome of coin tosses, black-and-white image pixels, binary representation of integers, etc.

- Goal: learn a model that predicts E_j from the spin configs.

Solving the problem with linear regression (lr):

- Without any prior knowledge of the dataset's origin, physics intuition may suggest: look for a spin model w/pairwise interactions between every pair of variables. That is we choose the following model class:

$$H_{model}[S^i] = - \sum_{j=1}^L \sum_{k=1}^L J_{j,k} S_j^i S_k^i \quad (79)$$

- Goal: Determine the interaction matrix $J_{j,k}$ by applying linear regression on the dataset \mathcal{D} = a well defined problem, since the unknown $J_{j,k}$ enters linearly into the definition of the Hamiltonian.
- To this end we cast the above ansatz into the more familiar linear regression form:

$$H_{model}[S^i] = \mathbf{X}^i \cdot \mathbf{J} \quad (80)$$

- Where: \mathbf{X}^i represent all two-body interactions $\{S_j^i S_k^i\}_{j,k=1}^L$, and i runs over the samples in the dataset.
- To complete analogy: can represent the dot product by a single index $p = \{j, k\}$, i. e. $\mathbf{X}^i \cdot \mathbf{J} = X_p^i J_p$.
- Note the regression model doesn't include the minus sign.

In the following we apply ordinary least squares (OLS), Ridge and Lasso regression to the problem, and compare.

- Fig 16 shows the R^2 (=coefficient of determination, a regression performance measure. Optimal performance: $R^2 = 1$) of the three regression models.

$$R^2 = 1 - \frac{\sum_{i=1}^n |y_i^{true} - y_i^{pred}|^2}{\sum_{i=1}^n |y_i^{true} - \frac{1}{n} \sum_{i=1}^n y_i^{pred}|^2} \quad (81)$$

- A few remarks:
 - The regularization parameter λ affects the Ridge and LASSo regressions at scales separated by a few orders of magnitude. Therefore, considered good practice to always check performance for the given model and data as a func of λ .

- While the OLS and Ridge regression test curves are monotonic, the LASSO test curve is not - suggesting an optimal LASSO regularization parameter is $\lambda \approx 10^{-2}$. At this sweet spot, the Ising interaction weights \mathbf{J} contains only nearest-neighbour terms (as did the model the data was generated from).
- Choosing between Ridge and LASSO in this case = similar to fixing gauge degrees of freedom.
 - Recall the uniform nearest-neighbour interaction strength $J_{j,k} = J$ which we used to generate the data, was set to unity, $J = 1$.
 - Moreover, $J_{j,k}$ was NOT defined to be symmetric (we only used the $J_{j,j+1}$, but never the $J_{j,j-1}$ elements).
- Fig 17 shows the matrix representation of the learned weights $J_{j,k}$. OSL and Ridge learn nearly symmetric weights $J \approx -0.5$. Not surprising, since it amounts to taking into account both the $J_{j,j+1}$ and $J_{j,j-1}$ terms, and the weights are distributed symmetrically between them.
- LASSO, on the other hand, tends to break this symmetry.
- Thus, we see how different regularization schemes can lead to learning equivalent models but in different "gauges".
- Any info we have about symmetry of the unknown model that generated the data should be reflected in the definition of the model and the choice of regularization.

6.5 Convexity of a regularizer

Mentioned previously the analytical LASSO solution can be found by invoking its convexity. We here provide an intro to **convexity theory**. Recall

- A set $C \subseteq \mathbb{R}^n$ is *convex* if for any $x, y \in C$ and $t \in [0, 1]$,

$$tx + (1 - t)y \in C \quad (82)$$

Aka, every line segment joining x, y lies entirely in C . A func $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is called convex if its domain, $\text{dom}(f)$, is a convex set, and for any $x, y \in \text{dom}(f)$ and $t \in [0, 1]$ we have

$$f(tx + (1 - t)y) \leq tf(x) + (1 - t)f(y) \quad (83)$$

That is, the func lies on or below the line segment joining its evaluation at x and y . This f is called *strictly convex* if this inequality holds strictly for $x \neq y$ and $t \in (0, 1)$.

- It turns out that *for convex functions, any local minimizer is a global minimizer*.
- Algorithmically, this means that in the optimization procedure, as long as we are "going down the hill" and agree to stop when we reach the minimum, then we have hit the global minimum.
- In addition to this, there is an abundance of rich theory regarding convex duality and optimality, which allow us to understand the solutions even before solving the problem itself.
- Let's examine the two regularizers Ridge and Lasso. They are both convex problems, but only Ridge is a strictly convex problem (assuming $\lambda > 0$).
- From convex theory this means we always have a unique sol for Ridge, but not necessarily for LASSO.
- It was recently shown that under mild conditions, the LASSO sol is indeed unique. Apart from this theoretical characterization, people have also introduced the notion of Elastic Net to retain the desirable properties of both LASSO and Ridge, which is now one of the standard tools for regression analysis and ML.

6.6 Bayesian formulation of linear regression

We formulate ls from a Bayesian point of view. We'll see that regularization in learning will emerge naturally as part of the Bayesian inference procedure.

From the linear regression setup, the data \mathcal{D} used to fit the regression model is generated through $y = \mathbf{x}^T \mathbf{w} + \epsilon$. We often assume ϵ is a Gaussian noise w/mean zero and variance σ^2 . To connect linear regression to the Bayesian framework, we often write the model as

$$p(y|\mathbf{x}, \boldsymbol{\theta}) = \mathcal{N}(y|\mu(\mathbf{x}), \sigma^2(\mathbf{x})) \quad (84)$$

Aka, our regression model is defined by a conditional probability that depends not only on data \mathbf{x} , but on some model parameters $\boldsymbol{\theta}$. Fex, if the mean a linear func of \mathbf{x} given by $\mu = \mathbf{x}^T \mathbf{w}$, and the variance is fixed $\sigma^2(\mathbf{x}) = \sigma^2$, then $\boldsymbol{\theta} = (\mathbf{w}, \sigma^2)$.

In statistics, many problems rely on estimation of some parameters of interest. Fex: suppose we're given the height data of 20 junior students from a regional high school, but we're interested in the average height of all high school juniors in the country. It's conceivable the data we're given are not representative of the student population as a whole. = necessary to devise a systematic way to perform reliable estimation. We here present the **Maximum Likelihood Estimation** (MLE), and show that least squares can be derived from this framework.

MLE is defined by

$$\hat{\boldsymbol{\theta}} = \operatorname{argmax}_{\boldsymbol{\theta}} \log p(\mathcal{D}|\boldsymbol{\theta}) \quad (85)$$

Using the assumption that samples are i.i.d., we can write the *log-likelihood* as

$$l(\boldsymbol{\theta}) \equiv \log p(\mathcal{D}|\boldsymbol{\theta}) = \sum_{i=1}^n \log p(y_i|\mathbf{x}_i, \boldsymbol{\theta}) \quad (86)$$

Inserting p as defined above we get

$$l(\boldsymbol{\theta}) = -\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \mathbf{x}_i^T \mathbf{w})^2 - \frac{n}{2} \log(2\pi\sigma^2) \quad (87)$$

$$= -\frac{1}{2\sigma^2} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2 + \text{const} \quad (88)$$

By comparing this to the expression for $\hat{\mathbf{w}}_{LS}$ it's clear that performing least squares is the same as maximizing the log-likelihood of this model.

What about adding regularization? We earlier introduced the *maximum a posteriori probability (MAP) estimate*. We now show it actually corresponds to regularized linear regression, where the choice of prior determines the regularization. Recall Baye's rule

$$p(\boldsymbol{\theta}|\mathcal{D}) \propto p(\mathcal{D}|\boldsymbol{\theta})p(\boldsymbol{\theta}) \quad (89)$$

Now instead of maximizing the log-likelihood $l(\boldsymbol{\theta}) = \log p(\mathcal{D}|\boldsymbol{\theta})$, let's maximize the log posterior $\log p(\boldsymbol{\theta}|\mathcal{D})$. Invoking Baye's rules, the MAP estimator becomes

$$\hat{\boldsymbol{\theta}}_{MAP} \equiv \operatorname{argmax}_{\boldsymbol{\theta}} \log p(\mathcal{D}|\boldsymbol{\theta}) + \log p(\boldsymbol{\theta}) \quad (90)$$

Review of Bayes terms:

- Likelihood function = $p(\mathbf{X}|\mathbf{w})$
- Prior distribution = $p(\mathbf{w})$
- Posterior distribution = $p(\mathbf{w}|\mathbf{X})$

Suppose we use the Gaussian prior (= the *conjugate prior* that gives a Gaussian posterior. For a given likelihood, conjugacy guarantees the preservation of prior distribution at the posterior level. Fex: for a Gaussian(Geometric) likelihood with a Gaussian(Beta) prior, the posterior distribution is still Gaussian(Beta) dist.) with zero mean and variance τ^2 , namely $p(\mathbf{w}) = \prod_j \mathcal{N}(w_j|0, \tau^2)$, we can recast the MAP estimator into (constant terms that don't depend on the parameters have been dropped)

$$\hat{\boldsymbol{\theta}}_{MAP} = \operatorname{argmax}_{\boldsymbol{\theta}} \left[-\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \mathbf{x}_i^T \mathbf{w})^2 - \frac{1}{2\tau^2} \sum_{j=1}^n w_j^2 \right] \quad (91)$$

$$= \operatorname{argmax}_{\boldsymbol{\theta}} \left[-\frac{1}{2\sigma^2} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2 - \frac{1}{2\tau^2} \|\mathbf{w}\|_2^2 \right] \quad (92)$$

The equivalence between MAP estimation with a Gaussian prior and Ridge regression is established by comparing this expression to the one for $\hat{\mathbf{w}}_{\text{Ridge}}(t)$ with $\equiv \sigma^2/\tau^2$. There's an analogous derivation for LASSO.

6.7 Recap and a general perspective on regularizers

In this section we

- Explored least square regression with and without regularization
- Motivated the need for regularization due to poor generalization, in particular the "high-dimensional limit" ($p \gg n$).
- Instead of showing the in-sample and out-of-sample errors explicitly, we conducted numerical experiments in Notebook 3 on the diabetes dataset and showed that regularization typically leads to better generalization.
- Due to the equivalence between the constrained and penalized form of regularized regression (in LASSO and Ridge, but not generally true in cases such as L^0 penalization), we can regard the regularized regression problem as an un-regularized problem but on a constrained set of parameters.
- Since the size of the allowed parameter space (e. g. $\mathbf{w} \in \mathbb{R}^p$ when un-regularized vs. $\mathbf{w} \in C \subset \mathbb{R}^p$ when regularized) is roughly a proxy for model complexity, solving the regularized problem is in effect solving the un-regularized problem with a smaller model complexity class = we're less likely to overfit.
- We also showed connection (using a regularization function \Leftrightarrow the use of priors in Bayesian inference). It gives us more intuition about why regularization implies we're less likely to overfit the data:
 - Say you're a physics student taking a lab class. Experiment goal is to measure the behavior of different pendula and use it that to predict the formula/model that determines the period of oscillation. You would record many things (inc the lenght and mass) to give yourself the best possible chance of determining the unknown relationship, the room temperature, air currents, table vibrations etc.
 - = You have a high dimensional dataset. Even higher - you're aware of time of day, that it's Wednesday, your friends in attendance, the location, etc. You didn't write it down because of strongly held prior beliefs that none of those things affect the physics taking place in the room.
 - What's serving you here is the intuition that probably only a few things matter in the physics of pendula. You're approaching the experiment with prior beliefs about how many features you'll need to pay attention to in order to predict what'll happen when you swing an unknown pendulum.

Point being that we live in a high-dimensional world of info and while we have have good intuition in the pendulum case/well-known problems, often in the ML field we cannot say with any confidence a priori *what* the small list of things to write down will be, but we can at least use regularization to help us enforce that the list not be too long so that we don't end up predicting that the pendulum's period depends on Bob having a cold on Wednesdays.

- In both LASSO and Ridge there is a parameter λ involved. In principle, this **hyperparameter** is usually predetermined = not part of the regression process. We saw our learning performance and solution depends strongly on λ = vital to choose it properly. As discussed, one approach is to assume an *uninformative prior* on the hyperparameters, $p(\lambda)$, and average the posterior over all choices of λ following this distribution. However, large computational cost. = simpler to choose the regularization parameter through some optimization procedure.
- Emphasize: linear regression can be applied to model non-linear relationship between input and response. Done by replacing the input \mathbf{x} w/some nonlinear func $\phi(\mathbf{x})$. Doing so preserves the linearity as a func of the parameters \mathbf{w} , since model is defined by the inner product $\phi^T(\mathbf{x})\mathbf{w}$. This method = *basis function expansion*.
- Recent years: surge of interest in understanding generalized linear regression models from a statistical physics perspective. Much has focused on understanding high-dimensional linear regression and compressed sensing. On a technical level, this research imports and extends the machinery of spin glass physics (replica method, cavity method, and message passing) to analyze high-dimensional linear models. A rich area of activity at the intersection of physics, computer science, information theory and ML.

7 Logistic regression

Have so far focused on continuous output. Now: classification: outcomes=discrete variables (i.e. categories). Fex, detect whether cat/dog in a picture. Or given a spin config, say the 2D Ising model, identify it's phase (e.g. ordered/disordered). Logistic regression deals with binary, dichotomous(=dividing into two parts) outcomes (e. g. true/false, success/failure, etc). Worth noting: logistic regression also commonly used in modern supervised deep learning models. Section structure:

- Define logistic regression and derive its corresponding cost function (the cross entropy) using a Bayesian approach and discuss its minimization.
- Generalize logistic regression to the case of multiple categories which is called *Softmax regression*.
- Demonstrate logistic regression via application to three different problems:
 - Classifying phases of the 2D Ising model

- The SUSY dataset
- MNIST handwritten digit classification

In this section, we consider the case where

- The dependent variables $y_i \in \mathbb{Z}$ are discrete and only take values from $m = 0, \dots, M - 1$ (i.e. M classes)
- Goal: predict the output classes from the design matrix $X \in \mathbb{R}^{n \times p}$ made of n samples, each of which bears p features. Of course primary goal = identify the classes to which new unseen samples belong.

Helpful to consider a slightly simple classifier before logistic regression:

- A linear classifier that categorizes examples using a weighted linear-combination of the features and an additive offset:

$$s_i = \mathbf{x}_i^T \mathbf{w} + b_0 \equiv \mathbf{x}_i^T \mathbf{w} \quad (93)$$

where $\mathbf{x}_i = (1, \mathbf{x}_i)$ and $\mathbf{w}_i = (b_0, \mathbf{w}_i)$.

- This func takes values on the entire real axis. With logistic regression however the labels y_i are discrete variables.
- One simple way to get a discrete output is to have sign functions that map the output of a linear regressor to $\{0, 1\}$, $f(s_i) = \text{sign}(s_i) = 1$ if $s_i \geq 0$ and 0 otherwise.
- = The "perceptron" in ML.

Extremely simple model, and it's favorable in many cases (e.g. noisy data) to have a "soft" classifier that outputs the probability of a given category:

- Fex: Given \mathbf{x}_i , the classifier outputs the probability of being in category m .
- One such func: The logistic (or **sigmoid**) function:

$$f(s) = \frac{1}{1 + e^{-s}} \quad (94)$$

Note that $1 - f(s) = f(-s)$ will be useful shortly.

7.1 The cross-entropy as a cost function for logistic regression

- Perceptron = example of a "hard classification": each datapoint deterministically assigned to a category.

- In many cases favorable to have a "soft" classifier that outputs the probability of a given category rather than a single value. Logistic regression is the most canonical example of a soft classifier.
- In logistic regression, the prob that a data point \mathbf{x}_i belongs to a category $y_i = \{0, 1\}$ is given by

$$P(y_i = 1|\mathbf{x}_i, \boldsymbol{\theta}) = \frac{1}{1 + e^{-\mathbf{x}_i^T \mathbf{w}}} \quad (95)$$

$$P(y_i = 0|\mathbf{x}_i, \boldsymbol{\theta}) = 1 - P(y_i = 1|\mathbf{x}_i, \boldsymbol{\theta}) \quad (96)$$

where $\boldsymbol{\theta} = \mathbf{w}$ are the weights we wish to learn from the data.

- For intuition:
 - Consider a collection of non-interacting two-state systems coupled to a thermal bath (e.g. a collection of atoms that can be in two states).
 - Denote the state of system i by a binary variable: $y_i \in \{0, 1\}$.
 - From elementary statistical mechanics, we know that if the two states have energies ϵ_0 and ϵ_1 the prob for finding the system in a state y_i is just:

$$P(y_i = 0) = \frac{e^{-\beta\epsilon_0}}{e^{-\beta\epsilon_0} + e^{-\beta\epsilon_1}} = \frac{1}{1 + e^{-\beta\Delta\epsilon}} \quad (97)$$

$$P(y_i = 1) = 1 - P(y_i = 0) \quad (98)$$

- Notice in these expressions, as is often the case in physics, only energy differences are observable.
 - If the difference in energies between two states is given by $\Delta\epsilon = \mathbf{x}_i^T \mathbf{w}$, we recover the expressions for logistic regression.
- We'll use this mapping between partition functions and classification to generalize the logistic regressor to soft-max regression later this section.
- Notice in terms of the logistic func, we can write

$$P(y_i = 1) = f(\mathbf{x}_i^T \mathbf{w}) = 1 - P(y_i = 0) \quad (99)$$

- We now define the cost func for logistic regression using Maximum Likelihood Estimation (MLE). In MLE we choose parameters to maximize the prob of seeing the observed data.
 - Consider a dataset $\mathcal{D} = \{(y_i, \mathbf{x}_i)\}$ w/binary labels $y_i \in \{0, 1\}$ where the data points are drawn independently.

- The likelihood of seeing the data under our model is

$$P(\mathcal{D}|\mathbf{w}) = \prod_{i=1}^n [f(\mathbf{x}_i^T \mathbf{w})]^{y_i} [1 - f(\mathbf{x}_i^T \mathbf{w})]^{1-y_i} \quad (100)$$

from which we compute the log-likelihood

$$l(\mathbf{w}) = \log P(\mathcal{D}|\mathbf{w}) \quad (101)$$

$$= \sum_{i=1}^n y_i \log f(\mathbf{x}_i^T \mathbf{w}) + (1 - y_i) \log [1 - f(\mathbf{x}_i^T \mathbf{w})] \quad (102)$$

- The maximum likelihood estimator = the set of parameters that maximize the log-likelihood:

$$\hat{\mathbf{w}} = \operatorname{argmax}_{\theta} \sum_{i=1}^n y_i \log f(\mathbf{x}_i^T \mathbf{w}) + (1 - y_i) \log [1 - f(\mathbf{x}_i^T \mathbf{w})] \quad (103)$$

- Since the cost (error) func = the negative log-likelihood, we have

$$\mathcal{C} = -l(\mathbf{w}) \quad (104)$$

$$= \sum_{i=1}^n -y_i \log f(\mathbf{x}_i^T \mathbf{w}) - (1 - y_i) \log [1 - f(\mathbf{x}_i^T \mathbf{w})] \quad (105)$$

This is known in statistics as the *cross-entropy*.

- Final note: Just as in linear regression we usually supplement the cross-entropy with additional regularization terms, usually L_1 and L_2 regularization.

7.2 Minimizing the cross entropy

Cross entropy = a convex function of the weights \mathbf{w} = any local minimizer is a global minimizer. Minimizing the cost func leads to the equation

$$\mathbf{0} = \nabla \mathcal{C}(\mathbf{w}) = \sum_{i=1}^n [f(\mathbf{x}_i^T \mathbf{w}) - y_i] \mathbf{x}_i \quad (106)$$

where used the logistic func identity $\partial_z f(z) = f(z)[1 - f(z)]$. This equation defines a transcendental equation for \mathbf{w} , the solution of which, unlike linear regression, cannot be written on a closed form. → Must use numerical methods such as introduced in the GD chapter to solve this optimization problem.

7.3 Examples of binary classification

7.3.1 Identifying the phases of the 2D Ising model

- Goal: show how to employ logistic regression to classify the states of the 2D Ising model according to their phase of matter.
- Hamiltonian for the classical Ising model:

$$H = -J \sum_{\langle ij \rangle} S_i S_j, \quad S_j \in \{\pm 1\} \quad (107)$$

where the lattice site indices i, j run over all nearest neighbors of a 2D square lattice, and J is an interaction energy scale. Adopt periodic boundary conditions.

- Onsager proved this model undergoes a phase transition in the thermodynamic limit from an ordered ferromagnet with all spins alligned to a disordered phase at the critical temperature $T_c/J = 2/\log(1 + \sqrt{2}) \approx 2.26$. For any finite system size, this critical point is expanded to a critical region around T_c .
- Can one train a statistical classifier to distinguish between the two phases of the Ising model? This could be used to locate the position of the critical point in more complicated models.
- Given an Ising state, we would like to classify whether it belongs to the ordered or disordered phase, without other info than the spin config itself. This categorical ML problem is well suited for logistic regression and will thus consist of recognizing whether a given state is ordered by looking at its bit configs.
- For the purposes of applying logistic regression, the 2D spin state will be flattened out to a 1D array, so won't be possible to learn info on the structure of the contiguous (=sharing a common border/touching) ordered 2D domains. Such info can be incorporated using deep CNNs, see later chapter.
- We consider a 2D Ising model on a 40×40 square lattice, and use Monte-Carlo (MC) sampling to prepare 10^4 states at every fixed temperature T out of a pre-defined set. We also assign a label to each state according to its phase: 0=disordered, 1=ordered.
- Near T_c the ferromagnetic correlation length diverges, leading to, amongst other things, critical slowing of the MC algo. Identifying phases could also be harder in this region perhaps. With this in mind, consider the following three types of states:

- Ordered ($T/J < 2.0$)
- Near-critical ($2.0 \leq T/J \leq 2.5$)
- Disordered ($T/J > 2.5$)

Use ordered and disordered states to train, then evaluate the performance on unseen ordered, disordered and near-critical states.

- We deploy the *liblinear* routine (= default for Scikit’s logistic regression) and stochastic gradient descent (SGD) to optimize the logistic regression cost func with L_2 regularization.
- We define accuracy of the classifier = the percentage of correctly classified data points. Comparing the accuracy on the training and test data we can study the degree of overfitting.
- See Fig 21: Notice the small degree of overfitting as suggested by the training and testing accuracy curves being close together.
- The liblinear minimizer outperforms SGD on the training and test data, but not on the near-critical data for certain values of the regularization strength λ .
- Similar to the linear regression examples, we find there exists a sweet spot for the SGD regularization strength λ that results in optimal performance of the logistic regressor, at about $\lambda \approx 10^{-1}$
- Does the difficulty of the phase recognition depend on the temp of the queried sample? For states in the near-critical temp region (see fig 20) it’s no longer easy for the human eye to distinguish between ordered and disordered \rightarrow interesting to compare the training and test accuracies to the accuracy of the near-critical state predictions (recall: model has not been trained on near-critical states). Indeed the liblinear accuracy is about 7% smaller for the critical states compared to the test data.
- Important to note all of Scikit’s logistic regression solvers have in-built regularizers. Crucial in order to prevent overfitting.

7.3.2 SUSY

- In high energy physics experiments we hope to discover new particles. Need to sift through events and classify as either a signal of a new process or particle, or as a background event from already understood Standard Model processes. We don’t know for sure what process has occurred - only know the final state particles. But, we can try to determine parts of phase space that will have a high percentage of signal events.

- Typically done by using a series of simple requirements on the kinematic quantities of the final state particles, for example having one or more leptons with large momenta that are transverse to the beam line (p_T).
- Instead, here we'll use logistic regression in an attempt to find the relative probability that an event is from a signal or background event. Rather than using the kinematic quantities of final state particles directly, we'll use the output of our logistic regression to define a part of phase space that is enriched in signal events.
- The dataset has been produced using Monte Carlo simulations to contain events with two leptons (electrons or muons). Each event has the value of 18 kinematic variables ("features"): first 8 = direct measurements of final state particles (in this case the p_T , pseudo-rapidity η , and azimuthal angle ϕ of two leptons in the event and the amount of missing transverse momentum (MET) together with its azimuthal angle), last 10 = functions of the first 8 - high level features derived by physicists to help discriminate between the two classes. These high level features can be viewed as the physicist's attempt to use non-linear functions to classify the events, having been developed with formidable theoretical effort. We'll later revisit this problem with the tools of Deep Learning.
- Since we don't know the true underlying process, so our goal in these types of analysis is to find regions enriched in signal events. If we find an excess of events above what is expected, we can have confidence they come from the signal type we're searching for. → The two metrics of interest are
 - The efficiency of signal selection
 - The background rejection achieved
- Often rather than thinking about just a single working point, performance is characterized by Receiver Operator Curves (or ROC curves). They plot signal efficiency (true-positives) vs. background rejection (true-negatives) as a function of some continuous variable such as a threshold. Here that variable will be the output signal probability of our logistic regression.
- Fig 22 shows examples of these outputs (Using L^2 regularization with a regularization parameter of 10^{-5}): x-axis=logistic regression model's output "probability of this being a true signal event", y-axis=frequency. Two figures: one where the actual event inputs are signal events (meaning the $p=1$ point on the x-axis counts true-positive cases, while the $p=0$ point counts false-negative cases (Type II error)) and one plot where the event inputs are background and not signal (meaning the

$p=1$ point on the x-axis counts false-positive cases (Type I error), while the $p=0$ point counts true-negative cases). → Some signal events even look background like (false-false), and some background events look signal like (false-true) = further reason to characterize performance in terms of ROC curves.

- Fig 23: Examples of ROC curves using L^2 regularization for many different regularization parameters using either TensorFlow (top) or SciKit learn (bottom) when using the full set of 18 input variables. Notice
 - Minimal overfitting, partly because such a large training dataset (4.5 million events). How do they see this from the ROC curves..??
 - More importantly, the underlying data we’re working with: each input variable is an important feature.
- Is there utility to this increased sophistication?
 - Recall, even to the learning algo, some signal events and background events look similar. Can illustrate this by a plot comparing the p_T spectrum of the leading and subleading leptons for both signal and background events. Fig 24 shows these two distributions. While *some* signal events are easily distinguished, many live in the same part of phase space as the background. This effect can also be seen in Fig 22.
 - How much discrimination power is obtained by simply putting different requirements on the input variables rather than using ML techniques? In order to compare this (called cut-based strategy in the HEP field) to regression, different ROC curves have been made for
 - * Logistic regression with just the simple kinematic variables,
 - * Logistic regression w/the full set of variables, and
 - * Just putting requirements on the leading lepton p_T
 - Fig 25 shows a clear performance benefit from using logistic regression.
 - Note in the cut-based approach we have only used one variable where we could have put requirements on all of them. While putting more requirements would increase background rejection, it would also decrease signal efficiency. → The cut-based approach will never yield as strong discrimination as logistic regression.
 - Also interesting in fig 25 - the higher-order variables noticeably help the ML techniques.

7.4 Softmax regression

- We generalize here from binary to multi-class classification.
- One approach is: treat the label as a vector $\mathbf{y}_i \in \mathbb{Z}_2^M$ = a binary bit string of length M . For $\mathbf{y}_i = (1, 0, \dots, 0)$ means the sample \mathbf{x}_i belongs to class 1.
- The prob of \mathbf{x}_i being in class m' is given by

$$P(y_{im'} = 1 | \mathbf{x}_i, \boldsymbol{\theta}) = \frac{e^{-\mathbf{x}_i^T \mathbf{w}_{m'}}}{\sum_{m=0}^{M-1} e^{-\mathbf{x}_i^T \mathbf{w}_m}} \quad (108)$$

where $y_{im'} \equiv [\mathbf{y}_i]_{m'}$ refers to the m' -th component of vector \mathbf{y}_i . This is the **softmax** function.

- Therefore, the likelihood of this M -class classifier is simply

$$P(\mathcal{D} | \{\mathbf{w}_k\}_{k=0}^{M-1}) = \prod_{i=1}^n \prod_{m=0}^{M-1} [P(y_{im} = 1 | \mathbf{x}_i, \mathbf{w}_m)]^{y_{im}} \times [1 - P(y_{im} = 1 | \mathbf{x}_i, \mathbf{w}_m)]^{1-y_{im}} \quad (109)$$

- From this we similarly define the cost func:

$$E(\mathbf{w}) = - \sum_{i=1}^n \sum_{m=0}^{M-1} y_{im} \log P(y_{im} = 1 | \mathbf{x}_i, \mathbf{w}_m) + (1 - y_{im}) \log(1 - P(y_{im} = 1 | \mathbf{x}_i, \mathbf{w}_m)) \quad (110)$$

As expected, for $M = 1$ we recover the cross entropy for logistic regression.

7.5 An example of SoftMax classification: MNIST digit classification

- The MNIST dataset: 70000 handwritten digits, each of which is laid out on a 28×28 -pixel grid.
- Every pixel assumes one of 256 grayscale values, interpolating between white and black.
- 10 categories for the digits 0 through 9 \rightarrow SoftMax regression with $M = 10$
- Fig 26 shows the learned weights \mathbf{w}_k where k corresponds to class labels (i.e. digits).
- We'll come back to SoftMax in chapter 9.

8 Combining models

- One of most powerful and widely applied ideas in modern ML: ensemble methods = combine predictions from multiple, often weak, statistical models. They undergird many of the winning entries in data science competitions such as Kaggle, especially on structured datasets (neural networks generally perform better than ensemble methods on unstructured data, images, and audio).
- Even in neural nets context, it's common to combine predictions from multiple neural nets to increase performance on tough image classification tasks.
- On one hand, idea of training multiple models, then using a weighted sum of their predictions, very natural. On the other hand, can imagine the ensemble predictions can be much worse than that of the individual models - especially when pooling reinforces weak but correlated deficiencies in each individual predictor. → Important to understand when we expect ensemble methods to work.
- To that end, will revisit the bias-variance tradeoff, generalize it to an ensemble of classifiers. Will show key to determine when ensemble methods work = the degree of correlation between the models in the ensemble.

8.1 Revisiting the bias-variance tradeoff for ensembles

- This tradeoff summarizes the fundamental tension in ML between (model complexity \Leftrightarrow amount of training data needed to fit it).
- Key property emerging from this analysis: the correlation between models that constitute the ensemble. Important for two distinct reasons:
 - Holding the ensemble size fixed, averaging the predictions of correlated models reduces the variance less than averaging uncorrelated models.
 - In some cases, correlations between models within an ensemble can result in an *increase* in bias, offsetting any potential reduction in variance gained from ensemble averaging. We'll discuss this in the context of bagging below. One of the most dramatic examples of increased bias from correlations = the catastrophic predictive failure of almost all derivative models used by Wall Street during 2008 financial crisis.

8.1.1 Bias-variance decomposition for ensembles

- We'll discuss bias-variance tradeoff in context of continuous predictions such as regression, but many ideas carry over to classification tasks.
- Review of the bias-variance tradeoff in context of a single model:
 - $\mathbf{X}_{\mathcal{L}} = \{(y_j, \mathbf{x}_j), j = 1 \dots N\}$ = data in dataset
 - $y = f(\mathbf{x}) + \epsilon$ = noisy model from which we assume the true data is generated from. ϵ = normally distributed w/mean zero and standard deviation σ_{ϵ} .
 - $\hat{g}_{\mathcal{L}}(\mathbf{x})$ = a predictor formed by a statistical procedure we assume we have, e.g. least squares regression. The predictor gives our model's prediction using a dataset \mathcal{L} .
 - This estimator is chosen by minimizing a cost func which, for the sake of correctness, we take to be the squared error

$$\mathcal{C}(\mathbf{X}, g(\mathbf{x})) = \sum_i (\mathbf{y}_i - \hat{g}_{\mathcal{L}}(\mathbf{x}_i))^2 \quad (111)$$

- \mathcal{L} = dataset drawn from some underlying distribution that describes the data.
- $\{\mathcal{L}_j\}$ = many different datasets of the same size as \mathcal{L} drawn from this distribution.
- The corresponding $\hat{g}_{\mathcal{L}_j}(\mathbf{x})$ will differ from each other due to stochastic effects arising from the sampling noise. \rightarrow Can view our estimator $\hat{g}_{\mathcal{L}}(\mathbf{x})$ as a random variable (technically functional) and define an expectation value $E_{\mathcal{L}}$ the usual way.
- $E_{\mathcal{L}}$ = expectation value computed by drawing infinitely many different datasets $\{\mathcal{L}_j\}$ of the same size, fitting the corresponding estimator, then averaging over the results.
- E_{ϵ} = the expectation value over different instances of the "noise" ϵ .
- Can then decompose the expected generalization error as

$$E_{\mathcal{L}, \epsilon}[\mathcal{C}(\mathbf{X}, g(\mathbf{x}))] = Bias^2 + Var + Noise \quad (112)$$

where the bias=the deviation of the expectation value of our estimator (=the asymptotic value of our estimator in the limit of infinite data) from the true value:

$$Bias^2 = \sum_i (f(\mathbf{x}_i) - E_{\mathcal{L}}[\hat{g}_{\mathcal{L}}(\mathbf{x}_i)])^2 \quad (113)$$

The variance = how much our estimator fluctuates due to finite-sample effects:

$$Var = \sum_i E_{\mathcal{L}}[(\hat{g}_{\mathcal{L}}(\mathbf{x}) - E_{\mathcal{L}}[\hat{g}_{\mathcal{L}}(\mathbf{x})])^2] \quad (114)$$

The noise term = part of the error due to intrinsic noise in the data generation process that no statistical estimator can overcome:

$$Noise = \sum \sigma_{\epsilon_i}^2 \quad (115)$$

Let's generalize this to ensembles of estimators:

- $\mathbf{X}_{\mathcal{L}}$ = a given dataset
- θ = given hyper-parameters that parametrize member of our ensemble
- We'll consider a procedure that deterministically generates a model $\hat{g}_{\mathcal{L}}(\mathbf{x}_i, \theta)$ given $\mathbf{X}_{\mathcal{L}}$ and θ .
- Assume θ includes some random parameters that introduce stochasticity into our ensemble (e.g. an initial condition for stochastic gradient descent or a random subset of features or data points used for training)
- We'll be concerned with the expected prediction error of the *aggregate ensemble predictor*

$$\hat{g}_{\mathcal{L}}^A(\mathbf{x}_i, \{\theta\}) = \frac{1}{M} \sum_{m=1}^M \hat{g}_{\mathcal{L}}(\mathbf{x}_i, \theta_m) \quad (116)$$

- For future reference we'll define the mean, variance, covariance (= the connected correlation function in the language of physics) and the normalized correlation coefficient wrt θ of the estimators in our ensemble as

$$\begin{aligned} E_{\theta}[\hat{g}_{\mathcal{L}}(\mathbf{x}, \theta)] &= \mu_{\mathcal{L}, \theta}(\mathbf{x}) \\ E_{\theta}[\hat{g}_{\mathcal{L}}(\mathbf{x}, \theta)^2] - E_{\theta}[\hat{g}_{\mathcal{L}}(\mathbf{x}, \theta)]^2 &= \sigma_{\mathcal{L}, \theta}^2(\mathbf{x}) \\ E_{\theta}[\hat{g}_{\mathcal{L}}(\mathbf{x}, \theta_m) \hat{g}_{\mathcal{L}}(\mathbf{x}, \theta_{m'})] - E_{\theta}[\hat{g}_{\mathcal{L}}(\mathbf{x}, \theta_m)]^2 &= C_{\mathcal{L}, \theta_m, \theta_{m'}}(\mathbf{x}) \\ \rho(\mathbf{x}) &= \frac{C_{\mathcal{L}, \theta_m, \theta_{m'}}(\mathbf{x})}{\sigma_{\mathcal{L}, \theta}^2} \end{aligned} \quad (117)$$

By definition we assume $m \neq m'$ in $C_{\mathcal{L}, \theta_m, \theta_{m'}}$.

We can now ask about the expected generalization (out-of-sample) error for the ensemble

$$E_{\mathcal{L}, \epsilon, \theta}[C(\mathbf{X}, \hat{g}_{\mathcal{L}}^A(\mathbf{x}))] = E_{\mathcal{L}, \epsilon, \theta}[\sum_i (\mathbf{y}_i - \hat{g}_{\mathcal{L}}^A(\mathbf{x}_i, \{\theta\}))^2] \quad (118)$$

As in the single estimator case, we decompose the error into a noise, bias and variance term. To see this, note that

$$E_{\mathcal{L},\epsilon,\theta}[\mathcal{C}(\mathbf{X}, \hat{g}_{\mathcal{L}}^A(\mathbf{x}))] = E_{\mathcal{L},\epsilon,\theta}[\sum_i (\mathbf{y}_i - f(\mathbf{x}_i) + f(\mathbf{x}_i) - \hat{g}_{\mathcal{L}}^A(\mathbf{x}_i, \{\theta\}))^2] \quad (119)$$

$$= \sum_i E_{\mathcal{L},\epsilon,\theta}[(\mathbf{y}_i - f(\mathbf{x}_i))^2 + (f(\mathbf{x}_i) - \hat{g}_{\mathcal{L}}^A(\mathbf{x}_i, \{\theta\}))^2 + 2(\mathbf{y}_i - f(\mathbf{x}_i))(f(\mathbf{x}_i) - \hat{g}_{\mathcal{L}}^A(\mathbf{x}_i, \{\theta\}))] \quad (120)$$

$$= \sum_i \sigma_{\epsilon_i}^2 + \sum_i E_{\mathcal{L},\theta}[(f(\mathbf{x}_i) - \hat{g}_{\mathcal{L}}^A(\mathbf{x}_i, \{\theta\}))^2] \quad (121)$$

where in the last line we used that $E_{\epsilon}[y_i] = f(\mathbf{x}_i)$ (why is only θ in the subindex of the E here?? How do the arithmetics of expec. values over several variables work? well probably the logical way, i.e. $E_{\mathcal{L},\epsilon,\theta} = (E_{\mathcal{L}} + E_{\epsilon} + E_{\theta})/3$ and probably ϵ was the only one they bothered to include here b.c. ϵ is the only variable we would expect $y = f(\mathbf{x}) + \epsilon$ to change over/be affected by) to eliminate the last term. Further decompose the second term as

$$E_{\mathcal{L},\theta}[(f(\mathbf{x}_i) - \hat{g}_{\mathcal{L}}^A(\mathbf{x}_i, \{\theta\}))^2] = E_{\mathcal{L},\theta}[(f(\mathbf{x}_i) - E_{\mathcal{L},\theta}[\hat{g}_{\mathcal{L}}^A(\mathbf{x}_i, \{\theta\})] + E_{\mathcal{L},\theta}[\hat{g}_{\mathcal{L}}^A(\mathbf{x}_i, \{\theta\})] - \hat{g}_{\mathcal{L}}^A(\mathbf{x}_i, \{\theta\}))^2] \quad (122)$$

$$= E_{\mathcal{L},\theta}[(f(\mathbf{x}_i) - E_{\mathcal{L},\theta}[\hat{g}_{\mathcal{L}}^A(\mathbf{x}_i, \{\theta\})])^2 + E_{\mathcal{L},\theta}[(E_{\mathcal{L},\theta}[\hat{g}_{\mathcal{L}}^A(\mathbf{x}_i, \{\theta\})] - \hat{g}_{\mathcal{L}}^A(\mathbf{x}_i, \{\theta\}))^2 + 2E_{\mathcal{L},\theta}[(f(\mathbf{x}_i) - E_{\mathcal{L},\theta}[\hat{g}_{\mathcal{L}}^A(\mathbf{x}_i, \{\theta\})])(E_{\mathcal{L},\theta}[\hat{g}_{\mathcal{L}}^A(\mathbf{x}_i, \{\theta\})] - \hat{g}_{\mathcal{L}}^A(\mathbf{x}_i, \{\theta\})))] \quad (123)$$

$$= (f(\mathbf{x}_i) - E_{\mathcal{L},\theta}[\hat{g}_{\mathcal{L}}^A(\mathbf{x}_i, \{\theta\})])^2 + E_{\mathcal{L},\theta}[(\hat{g}_{\mathcal{L}}^A(\mathbf{x}_i, \{\theta\}) - E_{\mathcal{L},\theta}[\hat{g}_{\mathcal{L}}^A(\mathbf{x}_i, \{\theta\})])^2] \quad (124)$$

$$= Bias^2(\mathbf{x}_i) + Var(\mathbf{x}_i) \quad (125)$$

where we defined the bias and variance of the aggregate predictor as

$$Bias^2(\mathbf{x}) \equiv (f(\mathbf{x}) - E_{\mathcal{L},\theta}[\hat{g}_{\mathcal{L}}^A(\mathbf{x}, \{\theta\})])^2 \quad (126)$$

$$Var(\mathbf{x}) \equiv E_{\mathcal{L},\theta}[(\hat{g}_{\mathcal{L}}^A(\mathbf{x}, \{\theta\}) - E_{\mathcal{L},\theta}[\hat{g}_{\mathcal{L}}^A(\mathbf{x}, \{\theta\})])^2] \quad (127)$$

So far, ensemble calc almost identical to that of a single estimator. But, since the aggregate estimator is a sum of estimators, its variance implicitly depends on the correlations between the individual estimators in the ensemble. Using the definition of the aggregate estimator $\hat{g}_{\mathcal{L}}^A$ and those of the mean, variance,

covariance and normalized correlation coefficients given earlier, we get

$$Var(\mathbf{x}) = E_{\mathcal{L},\theta}[(\hat{g}_{\mathcal{L}}^A(\mathbf{x}, \{\theta\}) - E_{\mathcal{L},\theta}[\hat{g}_{\mathcal{L}}^A(\mathbf{x}, \{\theta\})])^2] \quad (128)$$

$$= \frac{1}{M} \left[\sum_{m,m'} E_{\mathcal{L},\theta}[\hat{g}_{\mathcal{L}}^A(\mathbf{x}, \theta_m) \hat{g}_{\mathcal{L}}^A(\mathbf{x}, \theta_{m'})] - M^2 \sum_i [\mu_{\mathcal{L},\theta}(\mathbf{x})]^2 \right] \quad (129)$$

$$= \rho(\mathbf{x}) \sigma_{\mathcal{L},\theta}^2 + \frac{1 - \rho(\mathbf{x})}{M} \sigma_{\mathcal{L},\theta}^2 \quad (130)$$

This formula = key to understand the power of random ensembles. By using large ensembles ($M \rightarrow \infty$), can significantly reduce the variance, and for completely random ensembles where the models are uncorrelated ($\rho(\mathbf{x}) = 0$), maximally suppress the variance! \rightarrow Using the aggregate predictor beats down fluctuations due to finite-sample effects. Key = decorrelate the models as much as possible while still using a very large ensemble.

Can be worried this comes at the expense of a large bias. But no. When models in the ensemble completely random, the bias of the aggregate predictor = the expected bias of a single model:

$$Bias^2(\mathbf{x}) = (f(\mathbf{x}) - E_{\mathcal{L},\theta}[\hat{g}_{\mathcal{L}}^A(\mathbf{x}, \{\theta\})])^2 \quad (131)$$

$$= (f(\mathbf{x}) - \frac{1}{M} \sum_{m=1}^M E_{\mathcal{L},\theta}[\hat{g}_{\mathcal{L}}^A(\mathbf{x}, \theta_m)])^2 \quad (132)$$

$$= (f(\mathbf{x}) - \mu_{\mathcal{L},\theta})^2 \quad (133)$$

\rightarrow for a random ensemble one can always add more models without increasing the bias. This observation lies behind the immense power of random forest methods. For other methods inc bagging, we'll see that the bootstrapping procedure actually does increase the bias - but in many cases the increase is negligible compared to the reduction in variance.

8.1.2 Summarizing the theory and intuitions behind ensembles

Three distinct shortcomings fixed by ensemble methods: statistical, computational, representational.

- Statistical: When learning set too small, learning algo can typically find several models in the hypothesis space \mathcal{H} that all give the same performance on training data. Provided their predictions uncorrelated, averaging several models reduces risk of choosing wrong hypothesis.
- Computational: Many learning algos rely on some greedy assumption/local search that may get stuck in local optima. \rightarrow An ensemble made of individual models built from many different starting points may provide a better approximation of the true unknown function than any of the single models.

- **Representational:** In most cases, for a learning set of finite size, the true func cannot be represented by any of the candidate models in \mathcal{H} . By combining several models in an ensemble, may be possible to expand the space of representable funcs and better model the true func.

Combining models may come at the price of introducing more parameters to our learning procedure. But if the problem itself can never be learned through a simple hypothesis, then no reason to avoid applying a more complex model. Ensemble methods reduce the variance and often easier to train than a single complex model = powerful way of increasing representational power/expressivity.

Our analysis gives several intuitions for how we should construct ensembles.

- Should try to randomize ensemble construction as much as possible to reduce correlations between predictors in the ensemble. Ensures our variance reduced while minimizing increase in bias due to correlated errors.
- The ensembles will work best for procedures where the error of the predictor is dominated by the variance and not the bias. → These methods are especially well suited for unstable procedures whose results sensitive to small changes in the training data.

The ideas that using an ensemble allows us to reduce variance and that the procedure works best for unstable predictors which errors are dominated by variance caused by finite sampling rather than bias, may be carried over to classification tasks, even if this section has focused on continuous predictors such as regression.

8.2 Bagging

BAGGing = Bootstrap AGGregation = one of the simplest and most used ensemble methods. Imagine we have

- \mathcal{L} = a very large dataset that we could partition into M smaller datasets which we label $\{\mathcal{L}_1, \dots, \mathcal{L}_M\}$. If each partition is sufficiently large to learn a predictor, we can create an ensemble aggregate predictor, composed of predictors trained on each subset of the data.
- For continuous predictors like regression, this is just the average of all the individual predictors:

$$\hat{g}_{\mathcal{L}}^A = \frac{1}{M} \sum_{i=1}^M g_{\mathcal{L}_i}(\mathbf{x}) \quad (134)$$

- For classification tasks where each predictor predicts a class label $j \in \{1, \dots, J\}$, this is just the majority vote of all the predictors,

$$\hat{g}_{\mathcal{L}}^A(\mathbf{x}) = \operatorname{argmax}_j \sum_{i=1}^M I[g_{\mathcal{L}_i}(\mathbf{x}) = j] \quad (135)$$

where $I[g_{\mathcal{L}_i}(\mathbf{x}) = j]$ is an indicator func that is equal to one if $g_{\mathcal{L}_i}(\mathbf{x}) = j$ and zero otherwise. From our theoretical discussion we know this may significantly reduce the variance without increasing the bias.

- This form of aggregation clearly only works if enough data in each partition set \mathcal{L}_i . Consider the extreme limit where \mathcal{L}_i contains exactly one point. In this case, the base hypothesis $g_{\mathcal{L}_i}(\mathbf{x})$ (e.g. linear regressor) becomes extremely poor and the procedure above fails. One way to circumvent this shortcoming: resort to **empirical bootstrapping**, a resampling technique in statistics. Idea is to use sampling w/replacement to create new "bootstrapped" datasets $\{\mathcal{L}_1^{BS}, \dots, \mathcal{L}_M^{BS}\}$ from our original datasets \mathcal{L} . These bootstrapped datasets share many points, but due to the sampling w/replacement, are all somewhat different from each other.
- In the bagging procedure, we create an aggregate estimator by replacing the M independent datasets by the M bootstrapped estimators:

$$\hat{g}_{\mathcal{L}}^{BS}(\mathbf{x}) = \frac{1}{M} \sum_{i=1}^M g_{\mathcal{L}_i^{BS}}(\mathbf{x}) \quad (136)$$

and

$$\hat{g}_{\mathcal{L}}^{BS}(\mathbf{x}) = \operatorname{argmax}_j \sum_{i=1}^M I[g_{\mathcal{L}_i^{BS}}(\mathbf{x}) = j] \quad (137)$$

- This bootstrapping procedure allows us to construct an approximate ensemble and thus reduce variance. For unstable predictors, this can significantly improve performance.
- Price we pay for using bootstrapped training datasets is an increase in the bias of our bagged estimators: note that as the number of datasets $M \rightarrow \infty$, the expectation wrt the bootstrapped samples converges to the empirical distribution describing the training data set $p_{\mathcal{L}}(\mathbf{x})$ (e.g. a delta func at each datapoint in \mathcal{L}) = in general different from the true generative distribution for the data $p(\mathbf{x})$.
- Fig 29: Bagging with a perceptron (linear classifier) as the base classifier that constitutes the elements of the ensemble. See that although

each individual classifier in the ensemble performs poorly, bagging these estimators yields reasonably good performance. Raises questions like *why bagging works* and *how many bootstrap samples are needed*. As mentioned,

- Bagging is effective on "unstable" learning algos = small changes in the training set result in large changes in predictions.
- When the procedure is unstable, the prediction error is dominated by the variance and one can exploit the aggregation component of bagging to reduce the prediction error.
- In contrast, for a stable procedure the accuracy is limited by the bias introduced by using bootstrapped datasets.
- = There is an instability-stability transition point beyond which bagging stops improving our prediction.

Brief introduction to bootstrapping:

- $\mathcal{D} = \{X_1, \dots, X_n\}$ = a finite set of n data points given as training samples
- Our job: construct measures of confidence for our sample estimates (e.g. the confidence interval, or mean-squared error of sample median estimator).
- One then first samples n points **with replacement** from \mathcal{D} to get a new set $\mathcal{D}^{*(1)} = \{X_1^{*(1)}, \dots, X_n^{*(1)}\}$, called a **bootstrap sample**, which possibly contains repetitive elements.
- Then repeat the same procedure to get in total B such sets: $\mathcal{D}^{*(1)}, \dots, \mathcal{D}^{*(B)}$.
- Next step: use these B bootstrap sets to get the **bootstrap estimate** of the quantity of interest.
- Fex: let $M_n^{*(k)} = \text{Median}(\mathcal{D}^{*(k)})$ be the sample median of bootstrap data $\mathcal{D}^{*(k)}$.
- Then we can construct the variance of the distribution of bootstrap medians as:

$$\hat{Var}_B(M_n) = \frac{1}{B-1} \sum_{k=1}^B (M_n^{*(k)} - \bar{M}_n^*)^2 \quad (138)$$

where

$$\bar{M}_n^* = \frac{1}{B} \sum_{k=1}^B M_n^{*(k)} \quad (139)$$

is the mean of the median of all bootstrap samples.

- Specifically it has been shown that in the $n \rightarrow \infty$ limit, the distribution of the bootstrap estimate will be a Gaussian centered around $\hat{M}_n(\mathcal{D}) = \text{Median}(X_1, \dots, X_n)$ w/standard deviation $\propto 1/\sqrt{n}$.
- \rightarrow The bootstrap distribution $\hat{M}_n^* - \hat{M}_n$ approximates fairly well the sampling distribution $\hat{M}_n - M$ from which we obtain the training data \mathcal{D} .
- Note: M = the median based on which the true dist \mathcal{D} is generated. Aka, if we plot the histogram of $\{M_n^{*(k)}\}_{k=1}^B$, we'll see that in the large n limit it can be well fitted by a Gaussian which sharp peaks at $\hat{M}_n(\mathcal{D})$ and vanishing variance defined by \hat{Var}_B above.
- Fig 28: illustration.
- An interpretation of all the M s involved here:
 - M = the median based on which the true distribution/sample \mathcal{D} is generated
 - $M_n = ?$
 - $\hat{M}_n = \hat{M}_n(\mathcal{D}) = \text{Median}(X_1, \dots, X_n)$ = the median of the sample \mathcal{D}
 - $M_n^{*(k)} = \text{Median}(\mathcal{D}^{*(k)})$ = the sample median of bootstrap data $\mathcal{D}^{*(k)}$.
 - $\bar{M}_n^* = \frac{1}{B} \sum_{k=1}^B M_n^{*(k)}$ = the mean of the median of all bootstrap samples.
 - $\hat{M}_n^* = ?$
 - $\hat{M}_n - M$ = the sampling distribution from which we obtain the training data $\mathcal{D} = ?$
 - $\hat{M}_n^* - \hat{M}_n$ = bootstrap distribution $= \approx \hat{M}_n - M = ?$

8.3 Boosting

- In bagging, the contribution of all predictors weighted equally in the bagged (aggregate) predictor. But in principle a myriad ways to combine different predictors. Sometimes might prefer autocratic approach that emphasizes the best predictors, other times may be better to opt for more "democratic" ways as is done in bagging.
- In boosting an ensemble of weak classifiers $\{g_k(\mathbf{x})\}$ is combined into an aggregate, boosted classifier. But unlike bagging, each classifier is associated with a weight α_k indicating how much it contributes to the

aggregate classifier

$$g_A(\mathbf{x}) = \sum_{K=1}^M \alpha_K g_K(\mathbf{x}) \quad (140)$$

where $\sum_K \alpha_K = 1$.

- For already discussed reasons, boosting, like all ensemble methods, works best when combining simple, high-variance classifiers into a more complex whole.
- Here focus on "adaptive boosting" or AdaBoost. Basic idea: form the aggregate classifier in an iterative process.
- Important: at each iteration we reweight the error function to reinforce data points where the aggregate classifier performs poorly. In this way we can successively ensure our classifier has good performance over the whole dataset.
- Discussion of AdaBoost procedure in greater detail:

- $\mathcal{L} = \{(\mathbf{x}_i, y_i), i = 1, \dots, N\}$ a given dataset where $\mathbf{x}_i \in \mathcal{X}$ and $y_i \in \mathcal{Y} = \{+1, -1\}$
- Our objective: find an optimal hypothesis/classifier $g : \mathcal{X} \rightarrow \mathcal{Y}$ to classify the data.
- Let $\mathcal{H} = \{g : \mathcal{X} \rightarrow \mathcal{Y}\}$ be the family of classifiers available in our ensemble.
- In the AdaBoost we're concerned with the classifiers that perform somehow better than "tossing a fair coin". = For each classifier, the family \mathcal{H} can predict y_i correctly at least half the time.
- Construct the boosted classifier as follows:
 - * **Initialize** $w_{t=1}(\mathbf{x}_n) = 1/N, n = 1, \dots, N$
 - * **For** $t = 1, \dots, T$, **do**:
 - Select from \mathcal{H} a hypothesis g_t that minimizes the weighted error

$$\epsilon_t = \sum_{i=1: g_t(\mathbf{x}_i) \neq y_i}^N w_t(\mathbf{x}_i) \quad (141)$$

- Let $\alpha_t = \frac{1}{2} \ln \frac{1-\epsilon_t}{\epsilon_t}$, update the weight for each data \mathbf{x}_n by

$$w_{t+1}(\mathbf{x}_n) \leftarrow w_t(\mathbf{x}_n) \frac{e^{-\alpha_t y_n g_t(\mathbf{x}_n)}}{Z_t} \quad (142)$$

where $Z_t = \sum_{n=1}^N w_t(\mathbf{x}_n) e^{-\alpha_t y_n g_t(\mathbf{x}_n)}$ ensures all weights add up to unity.

- **Output** $g_A(\mathbf{x}) = \text{sign}(\sum_{t=1}^T \alpha_t g_t(\mathbf{x}))$

8.4 Random forests

- A random forest is composed of a family of (randomized) tree-based classifier decision trees (= high-variance, weak classifiers that can be easily randomized, and as such, are ideally suited for ensemble-based methods).
- A **decision tree** uses a series of questions to hierarchically partition the data. Each branch consists of a question that splits the data into smaller subsets (e.g. is some feature larger than a given number?), with leaves (end points) of the tree corresponding to the ultimate partitions of the data.
- Goal: construct trees such that the partitions are informative about the class label (see fig 30). More complex decision trees lead to finer partitions that give improved performance on the training set. But, this generally leads to overfitting, limiting the out-of-sample performance. → Almost all decision trees use some form of regularization (e.g. maximum depth of the tree) to control complexity and reduce overfitting.
- Decision trees have extremely high variance, and often extremely sensitive to many details of the training data. Not surprising since they're learned by partitioning the training data. → Individual trees are weak classifiers. But, these same properties make them ideal for incorporation in an ensemble method.
- To create ensemble of trees, must introduce a randomization procedure: Power of ensembles to reduce variance only manifests when randomness reduces correlations between the classifiers within the ensemble. Randomness usually introduced into random forests in one of three distinct ways.
 - First: Use bagging and simply "bag" the decision trees by training each decision tree on a different bootstrapped dataset. Strictly speaking this procedure does not constitute a random forest but rather **bagged decision trees**.
 - Second: Only use a different random subset of the features at each split in the tree. This "feature bagging" is the distinguishing characteristic of random forests. It reduces correlations between decision trees that can arise when only a few features are strongly predictive of the class label.
 - Finally: Extremized random forests (ERFs) combine ordinary and feature bagging with an extreme randomization procedure where splitting is done randomly instead of using optimality criteria.

Even though this reduces the predictive power of each individual decision tree, it still often improves the predictive power of the ensemble because it dramatically reduces correlations between members and prevents overfitting.

- Fig 31: Examples of the kind of decision surfaces found by decision trees, random forests, and Adaboost.

8.5 Gradient boosted trees and XGBoost

- Idea: Use intuition from boosting and gradient descent (in particular Newton's method) to construct ensembles of decision trees. As in boosting, the ensembles are created by iteratively adding new decision trees to the ensemble.
- In gradient boosted trees, a central role is played by a cost func measuring the performance of our ensemble. At each step, we compute the gradient of the cost func wrt the predicted value of the ensemble and add trees that move us in the dir of the gradient.
- Ofc, this requires clever way of mapping gradients to decision trees. We give a brief overview of how this is done within XGboost.
- Starting point: A clever parametrization of decision trees (dts). We here use notation where the dt make cont. predictions (regression trees), though this can easily be generalized to classification tasks.
- Parametrize a decision tree $g_j(\mathbf{x})$ with K leaves by two quantities:
 - $q(\mathbf{x})$ = a func that maps each data point to one of the leaves of the tree, $q : \mathbf{x} \in \mathbb{R}^d \rightarrow \{1, 2, \dots, K\}$
 - $\mathbf{w} \in \mathbb{R}^T$ that assigns a predicted value to each leaf.

In other words, the dt's prediction for the datapoint \mathbf{x}_i is simply: $q(\mathbf{x}_i) = w_{q(\mathbf{x}_i)}$

- In addition to a parametrization of dts, we also must specify a cost func which measures predictions. The prediction of our ensemble for a datapoint (y_i, \mathbf{x}_i) is given by

$$\hat{y}_i = g_A(\mathbf{x}_i) = \sum_{j=1}^M g_j(\mathbf{x}_i) \quad (143)$$

where $g_j(\mathbf{x}_i)$ is the prediction of the j -th dt on on datapoint \mathbf{x}_i , and M is the number of members in the ensemble.

As discussed in the context of random trees above, without regularization, dts tend to overfit the data by dividing it into smaller and smaller partitions. Thus, our cost func is generally composed of two terms:

- A term that measures the goodness of predictions on each data-point, $l_i(y_i, \hat{y}_i)$, which is assumed to be differentiable and convex
- And for each tree in the ensemble, a regularization term $\Omega(g_j)$ that does not depend on the data

$$\mathcal{C}(\mathbf{X}, g_A) = \sum_{i=1}^N l(y_i, \hat{y}_i) + \sum_{j=1}^M \Omega(g_j) \quad (144)$$

where index i runs over data points and j runs over dts in our ensemble. In XGBoost, the regularization func is chosen to be

$$\Omega(g) = \gamma T + \frac{\lambda}{2} \sum_{l=1}^T w_l^2 \quad (145)$$

with γ and λ regularization parameters that must be chosen appropriately. Notice this regularization penalizes both large weights on the leaves (similar to L^2 -regularization) and having large partitions w/many leaves.

- As in boosting, we form the ensemble iteratively. For this reason we define a family of predictors \hat{y}_t as

$$\hat{y}_i^{(t)} = \sum_{j=1}^t g_j(\mathbf{x}_i) = \hat{y}_i^{(t-1)} + g_t(\mathbf{x}_i) \quad (146)$$

Note: by definition $y_i^{(M)} = g_A(\mathbf{x}_i)$. The central idea is that for large t , each dt is a small perturbation to the predictor (of order $1/K$) and hence we can perform a Taylor expansion on our loss func to second order:

$$\mathcal{C}_t = \sum_{i=1}^N l(y_i, \hat{y}_i^{(t-1)} + g_t(\mathbf{x}_i)) + \Omega(g_t) \quad (147)$$

$$\approx \mathcal{C}_{t-1} + \Delta \mathcal{C}_t \quad (148)$$

with

$$\Delta \mathcal{C}_t = b_i l(y_i, \hat{y}_i^{(t-1)}) g_t(\mathbf{x}_i) + \frac{1}{2} a_i g_t(\mathbf{x}_i)^2 + \Omega(g_t) \quad (149)$$

where

$$a_i = \partial_{\hat{y}_i^{(t-1)}} l(y_i, \hat{y}_i^{(t-1)}) \quad (150)$$

$$b_i = \partial_{\hat{y}_i^{(t-1)}}^2 l(y_i, \hat{y}_i^{(t-1)}) \quad (151)$$

We then choose the t -th dt g_t to minimize $\Delta \mathcal{C}_t$. Almost identical to how we derived the Newton method update in the gradient descent section.

- Can actually derive an expression for the parameters of g_t that minimize $\Delta\mathcal{C}_t$ analytically. To simplify notation: define
 - The set of points \mathbf{x}_i that get mapped to leaf j : $I_j = \{i : q_t(\mathbf{x}_i) = j\}$
 - The functions $B_j = \sum_{i \in I_j} b_i$ and $A_j = \sum_{i \in I_j} a_i$.

In terms of these quantities we can write

$$\Delta\mathcal{C}_t = \sum_{j=1}^T [B_j w_j + \frac{1}{2}(A_j + \lambda_j) w_j^2] + \lambda T \quad (152)$$

where we made the t -dependence of all parameters implicit.

- To find optimal w_j , just as in Newton's method take gradient of the above expression wrt w_j and set this equal to zero, to get

$$w_j^{opt} = -\frac{B_j}{A_j + \lambda} \quad (153)$$

plugging back into $\Delta\mathcal{C}_t$ gives

$$\Delta\mathcal{C}_t^{opt} = -\frac{1}{2} \sum_{j=1}^K \frac{B_j^2}{A_j + \lambda} + \gamma T \quad (154)$$

→ Clear that $\Delta\mathcal{C}_t^{opt}$ measures the in-sample performance of g_t and we should find the dt that minimizes this value. Could in principle enumerate all possible trees over the data and find the tree that minimizes $\Delta\mathcal{C}_t^{opt}$.

- However, in practice, this is impossible. Instead, an approximate greedy algo is run that optimizes one level of the tree at a time by trying to find optimal splits of the data. Leads to a tree that is a good local minimum of $\Delta\mathcal{C}_t^{opt}$ which is then added to the ensemble.
- Emphasize this is only high level sketch of how the algo works. In practice, additional regularization such as shrinkage and feature subsampling is also used. In addition, there are many numerical and technical tricks used for the approximate algo and how to find splits of the data that give good dt s.

8.6 Application to the Ising model and Supersymmetry Datasets

- Now illustrate with using two physics examples (we previously analyzed both using logistic regression):

- Classifying the phases of the spin configs of the 2D-Ising model above and below the critical temp using random forests and
 - Classifying Monte-Carlo simulations of collision events in the SUSY dataset as supersymmetric or standard using an XGBoost implementation of gradient-boosted trees.
- We show that on the Ising dataset, the RFs perform significantly better than logistic regression models whereas gradient boosted trees seem to yield an accuracy of about 80%, comparable to published results.
- Ising: We assign a label to each state according to its phase: 0=disordered, 1=ordered. Divide the dataset into three categories according to the temp at which samples are drawn: ordered ($T/J < 2.0$), near-critical ($2.0 \leq T/J \leq 2.5$) and disordered ($T/J > 2.5$).
- Use the ordered and disordered states to train a random forest, evaluate our learned model on a test set of unseen ordered and disordered states. Also ask how well our RF can predict the phase of samples drawn in the region (=predict whether the temp of a critical sample is above or below the critical temp). Since our model never trained on critical region samples, prediction in this region is a test of the algo's ability to generalize to new regions in phase space.
- Result in Fig 32. Used two types of RFs: one where the ensemble consists of coarse dts w/few leaves, another w/fine dts w/many leaves. Extremely high accuracy on training and test sets (over 99%) for both.
- However, notice the RF consisting of coarse trees perform extremely poorly on samples from the critical region whereas RF w/fine trees classifies critical samples w/accuracy of nearly 85%.
- Interestingly, and unlike with logistic regression, this performance in the critical region requires almost no parameter tuning. This because, as discussed above, RFs are largely immune to overfitting problems even as the number of estimators in the ensemble becomes large. Increasing the number of estimators in the ensemble does increase performance - but at a large cost in computational time.
- SUSY: used the XGBoost implementation of gradient boosted trees to classify Monte-Carlo collisions from the SUSY dataset.
- With default parameters using a small subset of the data (100 000 out of the full 5 million samples), we achieved accuracy of 79%, which could be improved to nearly 80% after some fine tuning.
- Comparable to published results and those we obtained using logistic regression.

- One nice feature of ensemble methods such as XGBoost is that they automatically allow us to calc feature scores (Fscores) that rank the importance of various features for classification. Higher the Fscore, the more important the feature. Fig 33 shows feature scores for the production of electrically-charged supersymmetric particles (χ^\pm) which decay into W bosons and an electrically neutral supersymmetric particle χ^0 , which is invisible to the detector.
- The features are a mix of eight directly measurable quantities, as well as ten hand crafted features chosen using physics knowledge. Consistent with the physics of these supersymmetric decays in the lepton channel, we find the most informative features for classification are the missing transverse energy along the vector defined by the charged leptons (Axial MET) and the missing energy magnitude due to χ_0 .

9 An introduction to Feed-Forward Deep Neural Networks (DNNs)

- Long history, but remerged to prominence after a rebranding as "Deep Learning" in the mid 2000s.
- Truly caught attention in 2012 when Alex Krizhevsky, Ilya Sutskever and Geoff Hinton used a GPU-based DNN model (AlexNet) to lower the error rate on the ImageNet Large Scale Visual Recognition Challenge (ILSVRC) by an incredible twelve percent from 28% to 16%. Three years later an ML group from Microsoft achieved an error of 3.57% using an ultra-deep residual neural network (ResNet) with 152 layers. Since then, DNNs have become the workhorse technique for many image and speech recognition based ML tasks. Given rise to a number of high-level libraries and packages (Caffe, Kera, Pytorch, TensorFlow).
- Conceptually helpful to divide neural networks into four categories:
 - General purpose neural networks for supervised learning
 - Neural networks designed specifically for image processing, the most prominent example of this class being Convolutional Neural Networks (CNNs)
 - Neural networks for sequential data such as Recurrent Neural Networks (RNNs)
 - Neural networks for unsupervised learning such as Deep Boltzmann Machines.

In this chapter we'll discuss the first, and CNN and DPM in later chapters. While increasingly important for many applications, we omit RNNs.

- As with most intellectual fields experiencing rapid expansion, many commonly accepted heuristics turn out not to be as powerful as thought, and widely held beliefs not as universal as once imagined. Especially true in modern neural networks where results are largely empirical and heuristic and lack the firm footing of many earlier machine learning methods. Because of this, in this review have chosen to emphasize tried and true fundamentals, while pointing out what, from current vantage point, seem like promising new techniques.
- In physics, DNNs and CNNs have already found numerous applications:
 - Statistical physics: Detect phase transitions in 2D Ising and Potts models, lattice gauge theories, and different phases of polymers
 - Shown that DNNs can learn free-energy landscapes
 - Methods from statistical physics have been applied to the deep learning field to
 - * Study the thermodynamic efficiency of learning rules to explore the hypothesis space that DNNs span,
 - * Make analogies between training DNNs and spin glasses, and
 - * To characterize phase transitions w/respect to network topology in terms of errors.
 - In relativistic hydrodynamics, deep learning has been shown to capture features of non-linear evolution and has the potential to accelerate numerical simulations
 - In mechanics CNNs have been used to predict eigenvalues of photonic crystals
 - Recently, DNNs have been used to improve the efficiency of Monte-Carlo algorithms.
 - Deep learning has found interesting applications in quantum physics.
 - * Various quantum phase transitions can be detected and studied using DNNs and CNNs, including
 - The transverse-field Ising model
 - Topological phases
 - And even non-equilibrium many-body localization
 - * Representing quantum states as DNNs and quantum state tomography
 - * Study quantum and fault-tolerant error correction

- * Estimate rates of coherent and incoherent quantum processes
- * Recognition of state and charge configurations and auto-tuning in quantum dots
- * In quantum information theory, it's been shown one can perform gate-compositions with the help of neural nets
- * In lattice quantum chromodynamics, DNNs have been used to learn action parameters in regions of parameter space where PCA fails
- * Study of quantum control
- * Scattering theory to learn s-wave scattering length of potentials

9.1 Neural network basics

Neural networks/neural nets = neural-inspired nonlinear models for supervised learning. Natural, more powerful extensions of supervised learning methods such as linear and logistic regression and soft-max methods.

9.1.1 The basic building block: neurons

- The net's basic unit = a stylized "neuron" i that takes a vector of inputs $\mathbf{x} = (x_1, x_2, \dots, x_d)$ and produces a scalar output $a_i(\mathbf{x})$.
- A neural network consists of many such neurons stacked into layers, w/the output of one layer serving as the input of the next.
- First layer = input layer; middle layers= hidden layers; final layer = output layer.
- The exact function a_i varies depending on type of non-linearity used in the neural network. But, in essentially all cases a_i can be decomposed into a linear operation that weights the relative importance of the various inputs and a non-linear transformation $\sigma_i(z)$ which is usually the same for all neurons.
- The linear trans in almost all neural nets takes the form of a dot product with a set of neuron-specific weights $\mathbf{w}^{(i)} = (w_i^{(i)}, w_i^{(i)}, \dots, w_d^{(i)})$ followed by re-centering with a neuron-specific bias $b^{(i)}$

$$z^{(i)} = \mathbf{w}^i \cdot \mathbf{x} + b^{(i)} \quad (155)$$

- In terms of $z^{(i)}$ and the non-linear function $\sigma_i(z)$, we can write the full input-output function as

$$a_i(\mathbf{x}) = \sigma_i(z^{(i)}) \quad (156)$$

- Historically, common choices of nonlinearities included step-functions (perceptrons), sigmoids (i.e. Fermi functions), and the hyperbolic tangent.
- More recently, has become more common to use rectified linear functions (ReLU), leaky rectified linear units (leaky ReLU), and exponential linear units (ELUs).
- Different choices of non-linearities lead to different computational and training properties of neurons. Underlying reason is that we train neural nets using gradient descent based methods that require us to take derivatives of the neural input-output function with respect to the weights $\mathbf{w}^{(i)}$ and the bias $b^{(i)}$.
- Notice that the derivatives of the aforementioned non-linearities $\sigma(z)$ have very different properties. The derivative of the perceptron is zero everywhere except where the input is zero. \rightarrow This discontinuous behavior makes it impossible to train perceptrons using gradient descent.
- For this reason, until recently, the most popular choice of non-linearity was the tanh function or a sigmoid/Fermi function. But, this non-linearity choice has a major drawback. When the input weights become large, as they often do in training, the activation function **saturates** and the derivative of the output with respect to the weights tend to zero since $\frac{\partial \sigma}{\partial z} \rightarrow 0$ for $z \gg 1$.
- Such "vanishing gradients" are a feature of any saturating activation function (perceptron $\Theta(z)$, sigmoid $\frac{1}{1+e^{-z}}$, $\tanh \tanh(z)$), making it harder to train deep networks.
- In contrast, for a non-saturating activation function (such as ReLU $\max(0, z)$, leaky ReLU $(0.1z \text{ if } z \leq 0, z \text{ if } z \geq 0)$ or ELUs $(e^z - 1 \text{ if } z \leq 0, z \text{ if } z \geq 0)$), the gradients stay finite even for large inputs.

9.1.2 Layering neurons to build deep networks: network architecture

- Basic idea behind all neural networks: layer neurons in a hierarchical fashion, general structure of which is known as the network architecture.
- In the simplest feed-forward networks,
 - each neuron in the *input layer* of the neurons takes the inputs \mathbf{x} and produces an output $a_i(\mathbf{x} = \sigma_i(z^{(i)}))$ that depends on its current weights.
 - The outputs are then treated as inputs to the next *hidden layer*.

- This is usually repeated several times until one reaches the top or *output layer*. The output layer is almost always a simple classifier of the form discussed earlier: a logistic regression or soft-max function in the case of categorical data (i.e. discrete labels) or a linear regression layer in the case of continuous outputs.

Thus, the whole neural net can be thought of as a complicated nonlinear trans of the inputs \mathbf{x} into an output \hat{y} that depends on the weights and biases of all the neurons in the input, hidden and output layers.

- Hidden layers greatly expands the representational power of a neural net when compared to a simple soft-max or linear regression network.
- Perhaps the most formal expression of the increased representational power of neural networks (=expressivity) is the universal approximation theorem, stating: A neural network with a single hidden layer can approximate any continuous, multi-input/multi-output function with arbitrary accuracy. **The reader is strongly urged to read the beautiful graphical proof of the theorem in free online book (link).**
 - Basic idea of the proof: Hidden neurons allow neural networks to generate step functions with arbitrary offsets and heights. These can then be added together to approximate arbitrary functions.

The proof also makes clear that the more complicated a function, the more hidden units (and free parameters) are needed to approximate it. Hence, the applicability of the approximation theorem to practical situations should not be overemphasized.

- In physics, a good analogy are **matrix product states**, which can approximate any quantum many-body state to an arbitrary accuracy, provided the bond dimension can be increased arbitrarily - a severe requirement not met in any practical implementation of the theory.
- Modern neural networks generally contain multiple hidden layers (= the 'deep' in deep learning). There are many ideas why deep architectures favorable for learning.
 - Increasing number of layers = increases the number of parameters = increases the representational power.
 - Recent numerical experiments suggests that as long as the number of parameters is larger than the number of data points, certain classes of neural nets can fit arbitrarily labeled random noise samples. Suggests large neural nets of the kind used in practice can express highly complex functions.

- Adding hidden layers is also thought to allow neural nets to learn more complex features from the data. Work on CNNs suggests the first few layers learn simple, "low level" features that are then combined into higher-level, more abstract features in the deeper layers.
- Other works suggest it's computationally and algorithmically easier to train deep networks rather than shallow, wider nets, though this is still an area of major controversy and active research (**PAPER**).
- Choosing the exact network architecture remains an art that requires extensive numerical experimentation and intuition, and is often-times problem-specific. Both number of hidden layers and number of neurons in each layer can affect the performance of a neural network. Seems to be no single recipe for the right architecture for a neural net that works best. But, general rule of thumb that seems to be emerging: the number of parameters in the neural net should be large enough to prevent *underfitting* (**PAPER**).
- Empirically, the best architecture depends on
 - The task
 - The amount and type of data available
 - The computational resources at one's disposal.

Certain architectures are easier to train, while others might be better at capturing complicated dependencies in the data and learning relevant input features.

- Finally, have been numerous works that move beyond the simple deep, feed-forward neural nets discussed here. For example: Modern neural nets for image segmentation often incorporate "skip connections" that skip layers of the neural net, allowing information to directly propagate to a hidden or output layer, bypassing intermediate layers and often improving performance.

9.2 Training deep networks

- Basic procedure the same as we used for training simpler supervised learning algorithms, such as logistic and linear regression: Construct a cost/loss function and then use gradient descent to minimize the cost function. Neural nets differ from these simpler supervised procedures in that generally they contain multiple hidden layers that make taking the gradient more computationally difficult. Will return to this later when discussing the "backpropagation" algorithm for computing gradients.

- Like all supervised learning procedures, must first specify loss func.
 - Given a data point (\mathbf{x}_i, y_i) , the neural net makes a prediction $\hat{y}_i(\mathbf{w})$, where \mathbf{w} are the parameters of the neural network.
 - Recall that in most cases, the top output layer is either a continuous predictor or a classifier. Depending on this one must utilize a different loss function.
 - For continuous data:
 - * The loss funcs commonly used are identical to those in linear regression, include the mean squared error

$$E(\mathbf{w}) = \frac{1}{N} \sum_i (y_i - \hat{y}_i(\mathbf{w}))^2 \quad (157)$$

where N is the number of data points, and the mean absolute error (i.e. L_1 norm)

$$E(\mathbf{w}) = \frac{1}{N} \sum_i |y_i - \hat{y}_i(\mathbf{w})| \quad (158)$$

The full cost-func often includes additional terms that implement regularization (e.g. L_1 or L_2 regularizers).

- For categorical data:
 - * Most commonly loss-func is the cross-entropy, since the output layer is often taken to be a logistic classifier for binary data w/two types of labels, or a soft-max classifier if there are more than two types of labels. Cross-entropy already discussed extensively in logistic regression chapter.
 - * Recall that for binary data classification, the output of the top layer is the probability $\hat{y}_i(\mathbf{w}) = p(y_i = 1 | \mathbf{x}_i; \mathbf{w})$ that datapoint i is predicted to be in category 1.
 - * The cross-entropy between the true labels $y_i \in \{0, 1\}$ and the predictions is given by

$$E(\mathbf{w}) = - \sum_{i=1}^n y_i \log \hat{y}_i(\mathbf{w}) + (1 - y_i) \log [1 - \hat{y}_i(\mathbf{w})] \quad (159)$$

- * More generally, for categorical data, y can take on M values so that $y \in \{0, 1, \dots, M - 1\}$. For each datapoint i , define a vector y_{im} called a 'one-hot' vector, such that

$$y_{im} = \begin{cases} 1, & \text{if } y_i = m \\ 0, & \text{otherwise.} \end{cases}$$

Can also define the prob that the net assigns a datapoint to category m : $\hat{y}_{im}(\mathbf{w}) = p(y_i = m | \mathbf{x}_i; \mathbf{w})$. The categorical cross-entropy is then

$$E(\mathbf{w}) = - \sum_{i=1}^n \sum_{m=0}^{M-1} y_{im} \log \hat{y}_{im}(\mathbf{w}) + (1 - y_{im}) \log [1 - \hat{y}_{im}(\mathbf{w})] \quad (160)$$

As in linear and logistic regr, this loss func often supplemented by additional terms that implement regularization.

- Having identified architecture and cost func, time to train. Use gradient descent based (GD) methods to optimize cost func.
- Recall basic idea of GD = update the parameters \mathbf{w} to move in the dir of the gradient of the cost func, $\nabla_{\mathbf{w}} E(\mathbf{w})$. As discussed there are numerous variants. Most modern NN packages such as Keras allow user to specify which of these optimizers they would like to use in order to train the NN.
- Depending on the architecture, data, computational resources, different optimizers may work better, though vanilla SGD a good first choice.
- Unlike in linear and logistic regression, calc the gradients for a NN requires a special algo = backpropagation = backprop, forming the heart of the NN training.
- Backprop has been discovered multiple times independently, but was popularized for modern NNs in 1985.

9.3 High-level specification of a neural network using Keras

- Load the required packages
- Load the data (the MNIST digit data)
- Data preprocessing:
 - Split between training and testing, reshape
 - Standardize = normalize the input greyscale integer values to values between [0, 1].
 - Encode the classification labels using one-hot vectors, rather than integers (Keras provides function for doing this, "keras.utils.to_categorical")
- Build the network

- Create an instance of Keras's `Sequential()` class, and call it `model`. Allows us to build DNN's layer by layer. Use `add()` to attach layers to our `model`.
 - We focus on `Dense` layers for simplicity, but in subsequent examples, we'll see how to add dropout regularization and convolutional layers.
 - Every `Dense()` layer accepts as its first required argument an integer specifying the number of neurons.
 - Type of activation for the layer is defined using the `activation` optional argument, its input being the name of the activation function, for example `'relu'`, `'tanh'`, `'elu'`, `'sigmoid'`, `'softmax'`.
 - Must ensure number of input and output neurons for each layer match. → specify the shape of the input in the first layer of the model explicitly. The sequential construction of the model then allows Keras to infer the correct input/output dimensions of all hidden layers automatically. → we only need to specify the size of the softmax output layer to match the number of categories.
- Next, choose loss func to use for training. Chooses `categorical_crossentropy` defined in Keras' `losses` module (because we're dealing with classification i.e. cross-entropy, and because output data was cast in categorical form).
 - Choose SGD for optimization, available in Keras' `optimizers` module. Pass `lr` (learning rate) and `momentum` as optional arguments to the SGD func.
 - To test performance, look at a particular `metric` of performance. For instance, in categorical tasks typically looks at their `'accuracy'`, which is defined as the percentage of the correctly classified data points.
 - Use `compile()` to complete the model, with optional arguments for the `optimizer`, `loss` and the validation `metric`.
 - Training the DNN: a one-liner using the `fit()` method of the `Sequential` class. Two first required arguments = input and output data. Optional arguments: the `mini-batch_size`, the number of training `epochs`, and the test or `validation_data`. To monitor the training procedure for every epoch, we set `verbose=True`.

9.4 The backpropagation algorithm

Training requires us to calc the derivative of the cost func wrt all the parameters (the weight and biases of all neurons in all layers). A brute force calc

is **out of the question** since it requires us to calc as many gradients as parameters at each step of the gradient descent. Backprop = clever procedure that exploits layered structure of NNs to more efficiently compute gradients.

9.4.1 Deriving and implementing the backpropagation equations

- Backprop at its core: The ordinary chain rule for partial differentiation: if z is a function of n variables x_1, x_2, \dots, x_n and each of these variables are in turn functions of m variables t_1, t_2, \dots, t_m . Then for any variable t_i , $i = 1, 2, \dots, m$ we have that

$$\frac{\partial z}{\partial t_i} = \frac{\partial z}{\partial x_1} \frac{\partial x_1}{\partial t_i} + \frac{\partial z}{\partial x_2} \frac{\partial x_2}{\partial t_i} + \dots + \frac{\partial z}{\partial x_n} \frac{\partial x_n}{\partial t_i} \quad (161)$$

Backprop can be summarized using four equations.

- First: establish useful notation.
 - Assume there are L layers, with $l = 1, \dots, L$ indexing the layer.
 - Denote w_{jk}^l the weight for the connection from the k -th neuron in layer $l - 1$ to the j -th neuron in layer l .
 - Denote the bias of this neuron by b_j^l .
 - The activation a_j^l of the j -th neuron in the l -th layer can, by construction in a feed-forward NN, be related to the activities of the neurons in the layer $l - 1$ by the equation

$$a_j^l = \sigma\left(\sum_k w_{jk}^l a_k^{l-1} + b_j^l\right) = \sigma(z_j^l) \quad (162)$$

- We now define the four eqations making up backprop:
 1. By definition, the cost func E depends directly on the activities of the output layer a_j^L (which is our estimated response \hat{y} ?), but also indirectly on lower layer neuron activations, since a_j^L is calc from them. Define the error of neuron j in layer l , Δ_j^l , as the change in the cost func (why?) wrt the weighted input z_j^l :

$$\Delta_j^l = \frac{\partial E}{\partial z_j^l} = \frac{\partial E}{\partial a_j^l} \sigma'(z_j^l) \quad (163)$$

where $\sigma'(x)$ = the derivative of the non-linearity $\sigma(\cdot)$ wrt its input evaluated at x .

2. Notice the error func Δ_j^l can also be interpreted as the partial derivative of the cost func wrt the bias b_j^l , since

$$\Delta_j^l = \frac{\partial E}{\partial z_j^l} = \frac{\partial E}{\partial b_j^l} \frac{\partial b_j^l}{\partial z_j^l} = \frac{\partial E}{\partial b_j^l} \quad (164)$$

where we used that $\frac{\partial b_j^l}{\partial z_j^l} = 1$.

3. Derive the final two eq. using the chain rule. Since the error depends on neurons in layer l only through the activation of neurons in the subsequent layer $l + 1$, can use chain rule to write

$$\Delta_j^l = \frac{\partial E}{\partial z_j^l} \quad (165)$$

$$= \sum_k \frac{\partial E}{\partial z_k^{l+1}} \frac{\partial z_k^{l+1}}{\partial z_j^l} \quad (166)$$

$$= \sum_k \Delta_k^{l+1} \frac{\partial z_k^{l+1}}{\partial z_j^l} \quad (167)$$

$$= (\sum_k \Delta_k^{l+1} w_{kj}^{l+1}) \sigma'(z_j^l) \quad (168)$$

4. Derive the final eq by differentiating the cost func wrt w_{jk}^l

$$\frac{\partial E}{\partial w_{jk}^l} = \frac{\partial E}{\partial z_j^l} \frac{\partial z_j^l}{\partial w_{jk}^l} = \Delta_j^l a_k^{l-1} \quad (169)$$

- Together, these are the four backprop equations that relate the gradients of the activations a_j^l , the weighted inputs z_j^l and the errors Δ_j^l . They can be combined into a simple, computationally efficient algo to calc the gradient wrt all parameters: **The Backpropagation Algorithm**

1. **Activation at input layer:** Calc the activations a_j^1 of all the neurons in the input layer
2. **Feedforward:** starting w/the first layer, exploit the feed-forward architecture to compute z^l and a^l of each subsequent layer.
3. **Error at top layer:** calc the error of the top layer using equation 1.
4. **"Backpropagate" the error:** Use equation 3 to backpropagate the error backwards and calculate Δ_j^l for all layers.
5. **Calculate gradient:** Use eq. 2 and 4 to calc $\frac{\partial E}{\partial b_j^l}$ and $\frac{\partial E}{\partial w_{jk}^l}$.

- These basic ideas also underly almost all modern automatic differentiation packages such as Autograd (Pytorch).

9.4.2 Computing gradients in deep networks: what can go wrong with backprop?

- Even with backprop, gradient descent on large networks = extremely computationally expensive. The great advances in computational hardware (and the widespread use of GPUs) has made this a much less vexing problem than even a decade ago.

- Gradients can vanish and explode. *The problem of vanishing or exploding gradients.* Especially pronounced in NNs trying to capture long-range dependencies such as RNNs for sequential data. Can illustrate this by considering a simple network w/one neuron in each layer. Assume all weights equal, w . Behavior of the backprop for such a network can be inferred from repeatedly using eq. 3:

$$\Delta_j^1 = \Delta_j^L(w)^L \prod_{j=0}^{L-1} \sigma'(z_j) \quad (170)$$

Assume the magnitude $\sigma'(z_j)$ fairly constant, so we can approx $\sigma'(z_j) \approx \sigma'_0$. Then, notice that for large L , the error Δ_j^1 has very different behavior depending on the value of $w\sigma'_0$:

- $w\sigma'_0 > 1$: the errors of the gradient blow up.
- $w\sigma'_0 < 1$: the errors of the gradient vanish.
- $w\sigma'_0 \approx 1$ only now, and for neurons not saturated, will the gradient stay well behaved for deep networks.

This behavior holds true even in more complicated networks.

- Rather than considering a single weight, can ask about the eigenvalues (or singular values) of the weight matrices w_{jk}^l . In order for gradients to be finite for deep networks, we need these eigenvalues to stay near unity even after many gradient descent steps.
- In modern feedforward and ReLU NNs, this achieved by initializing the weights for the gradient descent in clever ways, and using non-linearities that don't saturate, such as ReLUs (recall: for saturating functions, $\sigma' \rightarrow 0$, causing the grad to vanish).
- Proper initialization and regularization schemes such as gradient clipping (cutting-off gradients w/very large values), and batch normalization also help mitigate the vanishing and exploding gradient problem.

9.5 Regularizing neural networks and other practical considerations

- DNNs, like all supervised learning algos, must navigate the bias-variance tradeoff.
- Regularization techniques play important role in ensuring DNNs generalize well to new data.
- Last five years seen a wealth of new specialized regularization techniques for DNNs beyond the simple L_1 and L_2 penalties, including

Dropout and Batch Normalization. Also, DNNs seem especially well-suited for implicit regularization that already takes place in the SGD. The implicit stochasticity and local-nature of SGD often prevents overfitting and spurious correlations in the training data, especially when techniques such as early stopping. We'll here give a brief overview over these techniques.

9.5.1 Implicit regularization using SGD: initialization, hyper-parameter tuning, and early stopping

- SGD = most commonly employed and effective optimizer for training NNs
- Acts as an implicit regularizer by introducing stochasticity that prevents overfitting.
- Important that weight initialization be chosen randomly, in order to break any leftover symmetries.
 - One common choice: drawing weights from a Gaussian centered around zero with some variance that scales inversely with number of inputs to the neuron. Since SGD a local procedure, as networks get deeper, *choosing a good weight initialization becomes increasingly important* to ensure the gradients are well behaved.
 - Choosing initialization with a variance too large or small causes gradients to vanish and the network to train poorly - even a factor of 2 can make huge difference in practice **PAPER**. → **important to experiment with different initialization variances**.
- Second important thing: appropriately choose the learning rate/step size by searching over five logarithmic grid points **PAPER**. If best performance occurs at edge of grid, repeat procedure until optimal learning rate is in the middle of the parameters.
- Common to center or whiten the input data (as we did for linear and logistic regression).
- Early stopping
 - Divide the training data into two portions, the dataset we train on and a smaller validation set that serves as a proxy for out-of-sample performance on the test set.
 - As we train model, we plot both training and validation error. We expect training error to continuously decrease during training. But validation error will eventually increase due to overfitting.

- Halt the training procedure when validation error starts to rise. This "early stopping" procedure ensures we stop training and avoid overfitting sample specific features in the data.

9.5.2 Dropout

- Basic idea: prevent overfitting by reducing spurious correlations between neurons by introducing a randomization procedure similar to that underlying ensemble models such as Bagging.
- Recall the basic idea behind ensemble methods: train an ensemble of methods that are created using a randomization procedure to ensure that the members of the ensemble are uncorrelated. This reduces the variance of statistical predictions without creating too much additional bias.
- Extremely costly to train an ensemble of NNs, both due to data needed and computational resources and parameter tuning.
- Dropout circumvents this, by randomly dropping out neurons (and their connections) from the NN during each step of the training.
- Typically, for each mini-batch in the grad descent step, a neuron is dropped from the NN with a probability p . The grad descent step then performed only on the weights of the "thinned" network of individual predictors.
- Since in the training, on average weights only present a fraction p of the time, predictions are made reweighing the weights by p : $\mathbf{w}_{test} = p\mathbf{w}_{train}$.
- The learned weights can be viewed as some "average" weight over all possible thinned NN. This averaging of weights is similar in spirit to the Bagging procedure discussed in the ensemble methods context.

9.5.3 Batch Normalization

- Basic inspiration: the long-known observation that training in NNs works best when the inputs are centered around zero wrt the bias. Because: it prevents neurons from saturating and gradients from vanishing in deep nets.
- In the absence of such centering, changes in parameters in lower layers can give rise to saturation effects in higher layers, and vanishing gradients.
- Idea: Introduce additional new "BatchNorm" layers that standardize the inputs by the mean and variance of the mini-batch.

- Consider a layer l with d neurons whose inputs are (z_1^l, \dots, z_d^l) . We standardize each dimension so that

$$\hat{z}_k^l = \frac{z_k^l - E[z_k^l]}{\sqrt{Var[z_k^l]}} \quad (171)$$

where the expectation and variance are taken over all samples in the mini-batch.

- One problem: this may change the representational power of the NN. Fex, for tanh non-linearities, it may force the network to live purely in the linear regime around $z = 0$. Since non-linearities are crucial to the representational power of DNNs, could dramatically alter its power.
- Thus, one introduces two new parameters γ_k^l and β_k^l for each neuron that can shift and scale the normalized input

$$\hat{y}_k^l = \gamma_k^l \hat{z}_k^l + \beta_k^l \quad (172)$$

These new parameters are then learned just like the weights and biases (just an extra layer for the backprop chain rule). We initialize the NN so at the beginning of training the inputs are standardized. Backprop then adjusts γ and β during training.

- In practice, this method considerably improves learning speed by preventing vanishing gradients. But, it also seems to serve as a powerful regularizer for reasons not fully understood. One plausible explanation: in this method, the gradient for a sample depends not only on the sample itself, but also on all the properties of the mini-batch. A single sample can occur in different mini-batches \rightarrow this introduces additional randomness into the training, which seems to help regularize training.

9.6 Deep neural networks in practice: examples

How to use NNs in practice.

9.6.1 Deep learning packages

- There are DNN packages in other languages than Python, fex Caffe in C++.
- Keras = high level framework, does not require any knowledge about inner workings of the underlying deep learning algos. But, for advanced applications, which may require more direct control over the operations in between layers, Keras' high level design may prove insufficient.

- Keras wraps the functionality of another package - TensorFlow.
- TensorFlow has become the preferred deep learning library. In it one constructs data flow graphs, where nodes=mathematical operations, edges=multidimensional tensors/data arrays. A DNN then thought of as a graph with a particular architecture. One needs to understand this concept well before one can truly unleash TensorFlow's full potential = steep learning curve and requires some time/perseverance.
- There are many other open source packages allowing for control over inter- and intra-layer operations, without need to introduce computational graphs. Fex: Pytorch. Offers libraries for automatic differentiation of tensors at GPU speed. As discussed, manipulating NNs boils down to fast array multiplication and contraction operations → the `torch.nn` library often does the job providing enough access and controllability to manipulate the linear algebra operations underlying DNNs.

9.6.2 Approaching the learning problem

- Typical procedure for using NNs:
 1. **Collect and pre-process the data.**
 2. **Define the model and its architecture.**
 3. **Choose the cost function and the optimizer.**
 4. **Train the model.**
 5. **Evaluate and *study* the model performance on the validation and test data.**
 6. **textbfAdjust the hyperparameters (and, if necessary, network architecture) to optimize performance for the specific dataset.**
 - Step 1: Getting data into appropriate form = inseparable part of the learning process.
 - One of the first questions = how to choose size of training vs test data.
 - * MNIST dataset, 10 classification categories: 80% training, 20% testing.
 - * ImageNet dataset, 100 categories: 50% training, 50% testing.
- Rule of thumb: more classification categories = closer should size of training and test data be, in order to prevent overfitting.
- * Once size of training data set, common to reserve 20% of it for validation, used for fine-tuning the hyperparameters of the model.

- How to choose the right hyperparameters to begin training with.
 - * Fex: according to Bengio, optimal learning rate = an order of magnitude lower than the smallest learning rate that blows up the loss.
 - * Also keep in mind, training the model can take a lot of time. This can severely slow down any progress on improving the model in Step 6. → usually a good idea to play with a small enough percentage of the training data to get a rough feeling about the correct hyperparameter regimes, the usefulness of the DNN architecture, and to debug one's code. Size of the "play set" should be such that training on it can be done fast and in real time to allow quickly adjusting the hyperparameters.
- Standardization of the dataset.
 - * Shown empirically that if the original values of the data differ by orders of magnitude, training can be slowed/impeded. Related to the vanishing/exploding gradient problem. → two tricks here: (i) All data should be mean-centered, i.e. from every data point we subtract the mean of the entire dataset (ii) Rescale the data (ensures the weights of the DNN are of a similar order of magnitude), for which there are two ways: if the data is approximately normally distributed, one can rescale by the standard deviation. Otherwise, it's typically rescaled by the max absolute value so the rescaled data lies in the interval $[-1, 1]$.
- Often, insufficient data serves as a major bottleneck on the ultimate performance of DNNs. Then one can consider data augmentation, i.e. distorting data samples from the existing dataset in some way to enhance size the dataset. If one knows how to do this, one already has partial information about the important features in the data.
- It's only when steps 1-5 are put together in step 6 that the real benefit of deep learning is revealed, compared to less sophisticated methods such as regression or bagging. The optimal choice of network architecture, cost func, and optimizer is determined by the properties of the training and test data, which are only revealed when we try to improve the model. A typical strategy for exploring the hyperparameter landscape is to use grid searches.
- This procedure can be applied to other ML tasks, not just DNNs. See chapter 11 for more useful hints on how to use the validation data.

9.6.3 SUSY dataset

- Look at the dataset we have previously studied with logistic regression and bagging. There is an interest in using deep learning methods to automatically discover collision features. Benchmark results using Bayesian Decision Trees from a standard physics package, and five-layer NNs using the dropout algorithm were presented in the original paper (reference) to compare the ability of deep learning to bypass the need of using such high-level features (what is meant by this sentence...?).
- Our goal: Study systematically the accuracy of a DNN classifier as a func of the
 - Learning rate
 - Dataset size

We here use Pytorch.

- We construct a DNN with two dense hidden layers of 200 and 100 neurons, respectively. We use ReLU activation between the input and hidden layers, and a softmax output layer. Apply dropout regularization on the weights of the DNN. Similar to MNIST, we use cross-entropy as cost func and minimize it using SGD with batches of size 10% of the training data size. We train the DNN over 10 epochs.
- Fig 40 shows accuracy of DNN on the test data as a func of the learning rate and the size of the dataset.
 - Considered good practice to start with a logarithmic scale to search through the hyperparameters, to get an overall idea for the order of magnitude of the optimal values.
 - In this example, performance peaks at the largest size of the dataset and a learning rate of 0.1, and is of the order of 80%.
 - For comparison, in the original study, the authors achieved $\approx 89\%$ by using the entire dataset with 5,000,000 points and a more sophisticated network architecture, trained using GPUs.

9.6.4 Phases of the 2D Ising model

- Study the problem of classifying the states of the 2D Ising model with a DNN, focusing on the model performance as a func of both the number of hidden neurons and the learning rate.
- We construct a minimalistic model for a DNN with a single hidden layer containing a number of hidden neurons. The network architecture thus includes a ReLU-activated input-layer, the hidden layer, and the

softmax output layer. We pick the categorical cross-entropy as cost func and minimize it using SGD with mini-batches of size 100. We train the DNN over 100 epochs.

- Fig 41 show the outcome of the grid search over a log-spaced learning rate and the number of neurons in the hidden layer. See that about 10 neurons are enough at learning rate of 0.1 to get a very high accuracy on the test set. But, if we aim at capturing the physics close to criticality, clearly more neurons are required to reliably learn the more complex correlations in the Ising states.

10 Convolutional Neural Networks (CNN)

- Core lesson in physics: exploit symmetries and invariances. Properties such as locality and translational invariance are often built directly into the physical laws. Our statistical physics model often directly incorporate everything we know about the physical system being analyzed. Fex: we know in many cases it's sufficient to consider only local couplings in our Hamiltonians, or work directly in momentum space if the system is translationally invariant.
- Like physical systems, many datasets and supervised learning tasks also possess additional symmetries and structure. Fex: consider a supervised learning task where we want to label images from some dataset as being pictures of cats or not. Our model must first learn features associated with cats. Because a cat a physical object, we know these features likely to be local (groups of neighboring pixels in the 2D image corresponding to whiskers, tails, eyes, etc.). Also know the cat can be anywhere in the image \rightarrow doesn't really matter where in the picture these features occur (though relative positions of features likely do matter). = a manifestation of translational invariance built into our supervised learning task.
- The all-to-all coupled NNs in the previous section fail to exploit this additional structure. Fex, consider image of the digit 'four' from the MNIST dataset. In the all-to-all coupled NNs used there, the 28×28 image was considered a 1D vector of size $28^2 = 796$. \rightarrow throws away lots of the spatial information.
- The NN community realized these problems, designed a class of NN architecture: CNNs, that take advantage of this additional structure (locality and translational invariance) **PAPER**.
- Interesting for physics: has been recently shown these CNN architectures are intimately related to models such as tensor networks **PAPER**

and, in particular, MERA-like architectures that are commonly used in physical models for quantum condensed matter systems **PAPER**.

10.1 The structure of convolutional neural networks

- CNN = translationally invariant NN that respects locality of the input data. Backbone of many modern deep learning applications.
- Reader encouraged to consult notes from the Stanford C231n CNN class this section has been based on
- There are two kinds of basic layers that make up a CNN:
 - A convolution layer that computes the convolution (a mathematical operations, see practical guide to image kernels) of the input with a bank of filters
 - Pooling layers that coarse-grain the input while maintaining locality and spatial structure.

For 2D data, a layer l is characterized by three numbers: height H_l , width W_l and depth D_l = the number of filters in that layer. All neurons corresponding to a particular filter have the same parameters (i.e. shared weights and bias).

- In general, will be concerned with local spatial filters (often called a receptive field in analogy with neuroscience) that takes as inputs a small spatial patch of the previous layer at all depths. Fex: a square filter of size F = a 3D array of size $F \times F \times D_{l-1}$. The convolution consists of running this filter over all locations in the spatial plane.
- To demonstrate:
 - Consider simple example consisting of a 1D input of depth 1 (fig 43). A filter of size $F \times 1 \times 1$ can be specified by a vector of weights w of length F .
 - The stride, S , encodes by how many neurons we translate the filter by when performing the convolution. In addition, it's common to pad the input with P zeros. For an input of width W , the number of neurons (outputs) in the layer is given by $(W - F + 2P)/S + 1$.
 - See link: visualization of the convolution procedure for a square input of unit depth.
 - After computing the filter, the output is passed through a non-linearity, a ReLU in fig 43. In practice, one often inserts a Batch-Norm layer before the non-linearity.

- These convolutional layers are interspersed with pooling layers that coarse-grain spatial information by performing a subsampling at each depth.
 - One common operation = max pool = the spatial dimensions are coarse grained by replacing a small region (say 2×2 neurons) by a single neuron whose output is the max value of the output in the region.
 - In **physics**, this pooling step very similar to the decimation step of RG. This generally reduces the output dimensions.
 - Fex: If region we pool over is 2×2 , then both the height and width of the output layer will be halved.
 - Generally, pooling do not reduce the depth of the convolutional layers because pooling is performed seperately at each depth.
 - Some studies suggests pooling might not be necessary, but it remains a staple of most CNNs.
- The convolutional and max-pool layers are generally followed by an all-to-all connected layer and a high level classifier such as soft-max.
 - Allows us to train CNNs as usual using the backprop algo. From a backprop perspective: CNNs almost identical to fully connected NN architectures except with tied parameters.
- Apart from introducing additional structure, such as translational invariance and locality, this conv structure also has important practical and computational benefits.
 - All neurons at given layer represent the same filter → can all be described by a single set of weights and biases. Reduces the number of free parameters by a factor of $H \times W$ at each layer.
 - Fex: For a layer with $D = 10^2$ and $H = W = 10^2$, this gives a reduction in parameters of nearly 10^6 !
 - → Can train much larger models than would otherwise be possible with fully connected layers.
 - Similar phenomena in physics: e.g. in translationally invariant systems we can parametrize all eigenmodes by specifying their momentum (wave number) and functional form (sin, cos, etc), while without translation invariance much more information is required.

10.2 Example: CNNs for the 2D Ising model

- Inclusion of spatial structure = important feature that can be exploited when designing NNs for studying physical systems.

- Used Pytorch to implement a simple CNN composed of
 - A single conv layer followed by a soft-max layer.
 - We varied the depth of the CNN layer from unity - a single set of weights and one bias - to a depth of 50 distinct weights and biases.
 - Trained using SGD for five epochs using a training set consisting of samples from far in the paramagnetic and ordered phases.
- Fig 45 results: T
 - The CNN achieved a 100% accuracy on the test set for all architectures, even for CNN of depth one.
 - Also checked the performance on samples drawn from the near-critical region of temperatures T slightly above and below the critical temperature T_c . The CNN performed admirably even on these critical samples with accuracy between 80% and 90%.
 - As with all ML and NNs, performance on parts of data missing from the training data considerably worse than on test data similar to training data. Highlights the importance of properly constructing an accurate training dataset and the considerable obstacles of generalizing to novel situations.
- Regarding the SUSY dataset, we stress that the absence of spatial locality in the collision features renders applying CNNs to that problem inadequate. (Do they mean that this particular provided input data don't include spatial info, or that spatial info is actually not part of a collision? I'm assuming the first.)

10.3 Pre-trained CNNs and transfer learning

- The immense success of CNNs for image recognition → training of huge networks on enormous datasets often by large industrial research teams from Google, Microsoft, Amazon etc. Many of these models known by name: AlexNet, GoogLeNet, ResNet, InceptionNet, VGGNet, etc.
- The trained models have been released, now available in standard packages such as the Torch Vision library in Pytorch, or the Caffe framework. They can be used directly as a basis for fine-tuning in different supervised image recognition tasks through a process called transfer learning.
- Transfer learning basic idea: The filters (receptive fields) learned by the conv layers of these networks should be informative for most image recognition based tasks, not just the ones they were originally trained

for. This turns out to be true in practice for many tasks one might be interested in.

- Three distinct ways one can take a pretrained CNN and repurpose it for a new task:
 - **Use CNN as fixed feature detector at top layer.** If new dataset we want to train on small and similar to original dataset, can simply use the CNN as a fixed feature detector and retrain our classifier = remove the classifier (soft-max) layer at the top of the CNN and replace it with a new classifier (linear SVM or soft-max) relevant to our problem. Here, the CNN serves as a fixed map from images to relevant features (the outputs of the top fully-connected layer right before the original classifier). This prevents overfitting on small, similar datasets. Often useful starting point for transfer learning.
 - **Use CNN as fixed feature detector at intermediate layer.** If dataset small and different from the one used to train the original model, the features at the top level might not be suitable for our dataset. → may want to instead use features in the middle of the CNN to train our new classifier. These features thought to be less fine-tuned and more universal (e.g. edge detectors). This is motivated by the idea that CNNs learn increasingly complex features the deeper one goes in the network (see discussion on representational learning in next section)
 - **Fine-tune the CNN.** If dataset large: in addition to replacing and training the classifier in the top layer, can also fine-tune the weights of the original CNN using backprop. May choose to freeze some of the weights in the CNN during the procedure or retrain all of them simultaneously.
- All this can be carried out easily using packages such as Caffe or the Torch Vision library in PyTorch.

11 High-level concepts in Deep Neural Networks

11.1 Organizing deep learning workflows using the bias-variance tradeoff

- Imagine you're given some data and asked to design an NN for learning how to perform a supervised learning task. What are the best practices for organizing a systematic workflow that allows us to efficiently do this?

- Here we present a simple deep learning workflow inspired by thinking about the bias-variance tradeoff (fig 46). This section draws heavily on tutorial from the Deep Learning School.
- **First thing** we would like to do: divide data in three: training set, validation/development/dev set, and test set. Use validation error as proxy for the test error in order to make tweaks to our model. OBS: do not use any test data to train the algo. This is a cardinal sin in ML.
- **Estimate optimal error rate (Bayes rate)**
 - Establish the difficulty of the task and the best performance one can expect to achieve. No algo can do better than the "signal" in the dataset. Fex: it's likely much easier to classify objects in high resolution images than in blurry, low-resolution ones. → Must establish a proxy or baseline for the optimal performance that can be expected from any algo.
 - In Bayesian statistics, this known as the Bayes rate. Since we don't know this *a priori*, we must get an estimate. For many tasks inc speech and object recognition, can approximate this by humans' performance on it. For more specialized task, we would like to ask how well experts, trained at the task, perform. This expert formance then serves as proxy for our Bayes rate.
- **Minimize underfitting (bias) on training data set.**
 - Having the Bayes rate, want to ensure we are using a sufficiently complex model to avoid underfitting on the training data.
 - In practice this means: comparing the training error to the Bayes rate. Since the training error doesn't care about generalization (variance), our model should approach the Bayes rate in the training set.
 - If not, bias of the DNN is too large and should try training the model longer and/or use larger model.
 - If none of these techniques work, likely the model architecture not well suited to the data, should modify the architecture in some way to better reflect the underlying structure of the data (symmetries, locality, etc.)
- **Make sure you are not overfitting.** Next, run our algo on the validation/dev set.
 - If error similar to training error rate and Bayes rate, we're done.

- If not, we are overfitting the training data. Possible solutions include: regularization and, importantly, collecting more data.
- If none of these work: likely has to change the DNN architecture.
- The result
 - If validation and test set drawn from same distributions, good performance on validation set should \Leftrightarrow good performance on the test set. (But typically slightly worse because the hyperparameters were fit to the validation data).
 - But, sometimes training and test data differ in subtle ways. Fex: they're collected using slightly different methods, or it's cheaper to collect data one way vs another. (But why on earth would you then divide training/test data based on this difference? Wouldn't you randomize it?) Rectification: Make two validation/dev sets, one from training data and one from test data. The difference in performance on the two quantifies the train-test mismatch. This can serve as another important diagnostic when using DNNs.

11.2 Why neural networks are so successful: three high-level perspectives on neural networks

As pointed out earlier, the field is rapidly expanding, and many of these perspective may out to be only partially true or even false. Nonetheless, included as guidepost for readers.

11.2.1 Neural networks as representation learning

- Powerful aspect of deep learning: The ability to learn relevant features with relatively little domain knowledge/minimal hand-crafting. Power of deep learning often stems from its ability to act like a black box: take in a large stream of data - find good features capturing the properties of the data we're interested in.
- This ability to learn good representations with very little hand-tuning = one of the most attractive properties of DNNs.
- Many of the other supervised-learning algos discussed (regression-based models, ensemble methods inc random forests or gradient-boosted trees) perform comparably or even better than NNs - but when using hand-crafted features with small-to-intermediate sized datasets.
- The hierarchical structure of deep learning = thought to be crucial to their ability to represent complex, abstract features. Fex: how analysis of CNNs for image classification suggests the lower-levels of the net learn elementary features, such as edge detectors, which are

then combined into higher levels of the net into more abstract, higher-level features (e.g. the famous example of a neuron that "learned to respond to cats").

- Has been shown more recently that CNNs can be thought of as **performing tensor decompositions on the data similar to those commonly used in numerical methods in modern quantum condensed matter. PAPER.**
- Interesting consequence of this thinking: One can train a CNN on one large dataset and the features it learn should also be useful for other supervised tasks → the ability to learn important and salient features directly from the data, then transfer the knowledge to new task. This ability to learn important, higher-level, coarse grained features is reminiscent of **ideas like the renormalization group (RG) in physics where the RG flows separate out relevant and irrelevant directions, and certain unsupervised deep learning architectures have a natural interpretation in terms of variational RG schemes. PAPER.**

11.2.2 Neural networks can exploit large amounts of data

- Data explosion
- DNNs are able to exploit the additional signal in large datasets for difficult supervised learning tasks.
- Fundamentally, modern DNNs are unique in that they contain millions of parameters, yet can still be trained on existing hardware. Complexity of DNNs (in terms of parameters) combined with their simple architecture (layer-wise connections) hit a sweet spot between expressivity (ability to represent very complicated functions) and trainability (ability to learn millions of parameters).
- Indeed, DNN's ability to exploit large datasets thought to differ from many other commonly employed supervised learning methods, like Support Vector Machines (SVMs). Fig 47: schematic depicting expected performance of DNNs of different sizes with the number of data samples and compares them to supervised learning algorithms such as SVMs or ensemble methods.
 - When amount of data small: DNNs offer no substantial benefit over these other methods and often performs worse.
 - But: large DNNs seem to be able to exploit additional data in a way other methods cannot.

- Fact that one doesn't have to hand engineer features makes the DNN even more well suited for handling large datasets.
- Recent theoretical results suggest that as long as a DNN large enough, should generalize well and not overfit **ARTICLE**.

11.2.3 Neural networks scale up well computationally

- Modern NNs can harness the immense computational capability that has occurred over the last few decades. Architecture of NNs naturally lends itself to parallelization and the exploitation of fast but specialized processors such as graphical processing units (GPUs).
- Google and NVIDIA set on a course to develop TPUs (tensor processing units) which will be specifically designed for the mathematical operations underlying deep learning architectures.
- The layered architecture of NNs also makes it easy to use modern techniques such as automatic differentiation that make it easy to quickly deploy them.
- Algos such as SGD and use of mini-batches make it easy to parallelize code and train much larger DNNs than was thought possible fifteen years ago.
- Many of these computational gains are quickly incorporated into modern packages with industrial resources → makes it easy to perform numerical experiments on large datasets, leading to further engineering gains.

11.3 Limitations of supervised learning with deep networks

Supervised learning using NNs has important limitations, like all statistical methods. Especially important when seeking to apply these methods to physics. Often, the same or better performance on a task can be achieved by using a few hand-engineered features (or even a collection of random features). Especially important for hard physics problems where data/Monte-Carlo samples maybe hard to come by. Some important limitations:

- **Need labeled data.** Like all supervised learning methods. Can be harder to acquire than unlabeled data (e.g. must pay human experts to label images).
- **Supervised neural networks are extremely data intensive.** Utility of DNNs extremely limited if data is hard to acquire or the datasets are small (hundreds to a few thousands samples). In this case, performance of other methods that utilize hand-engineered features can exceed that of DNNs.

- **Homogeneous data.** Almost all DNNs deal with homogeneous data of one type. Very hard to design architectures that mix and match data types (i.e. some continuous variables, some discrete variables, some time series). In applications beyond images, video and language, this is often what is required. In contrast, ensemble methods like random forests or gradient-boosted trees have no difficulty handling mixed data types.
- **Many physics problems are not about prediction.** Often not interested in solving prediction tasks such as classification. Want to learn something about the underlying distribution that generates the data. Then, often difficult to cast these ideas in a supervised learning setting. While the problems are related, it's possible to make good predictions with a "wrong" model. The model might or might not be useful for understanding physics.

Some of these remarks particular to DNNs, other shared by all supervised learning methods. Motivates the use of unsupervised methods which in part circumnavigate these problems.

12 Dimensional reduction and data visualization

- Will begin our foray into unsupervised learning by way of data visualization = an important tool in ML to identify structures such as
 - Correlations
 - Invariances (symmetries)
 - Irrelevant features (noise)
 in raw or processed data.
- Conceivably, capturing these properties could help us design better predictive models. In practice, however, the data we deal with is often high-dimensional = its visualization is impossible - daunting at best.
- Part of the complication = low-dimensional representation of high-dimensional data necessarily incurs information lost.
- Simple way to visualize data: pair-wise correlations (= pairwise scatter plots for all features). Useful in highlighting the important correlations between features when the number of features we are measuring is relatively small.
- In practice, we often have to perform *dimensional reduction* = project the data onto a lower dimensional space = *the latent space*.

- We discuss both linear and non-linear methods for dimensional reduction with applications in data visualization. The techniques can be used in many other applications also, inc lossy data compression and feature extraction.

12.1 Some of the challenges of high-dimensional data

- *High-dimensional data lives near the edge of sample space* Geometry in high-dimensional space can be counterintuitive. Example pertinent to ML:
 - Consider data distributed uniformly at random in a D -dimensional hypercube $\mathcal{C} = [-e/2, e/2]^D$, where e = the edge length.
 - Consider also a D -dimensional hypersphere \mathcal{S} of radius $e/2$ centered at the origin and contained within \mathcal{C} .
 - The prob that a data point \mathbf{x} drawn uniformly at random in \mathcal{C} is contained within \mathcal{S} is well approximated by the ratio of the volume of \mathcal{S} to that of \mathcal{C} : $p(\|\mathbf{x}\|_2 < e/2) \sim (1/2)^D$.
 - Thus, as the dimension of the feature space D increases, p goes to zero exponentially fast. = most of the data will concentrate outside the hypersphere, in the corners of the hypercube.
 - In **physics**, this basic observation underlies many properties of ideal gas such as the Maxwell distribution and the equipartition theorem.
- *Real-world data vs. uniform distribution* Fortunately, real-world data is not random or uniformly distributed!
 - Real data usually lives in a much lower dimensional space than the original space in which the features are being measured (see very helpful fig 48).
 - "The blessing of non-uniformity" (vs the curse of dimensionality)
 - Data will typically be locally smooth (= a local variation of the data will not incur a change in the target variable)
 - The idea = similar to statistical physics, where properties of most systems w/many degrees of freedom can often be characterized by low-dimensional 'order parameters'.
 - In thermodynamics, bulk properties of a gas of weakly interacting particles can be simply described by the thermodynamic variables that enter equation of states rather than astronomically large dynamical variables (= position and momentum) of each particle in the gas is another instantiation of this idea.
- *The crowding problem*

- When performing dimensional reduction, a common goal is to preserve pairwise distances between the data points from the original space to the latent space.
- Can be fairly well achieved if the *intrinsic* dimensionality of the data (=that in the original space) is the same as the dimension of the latent space. Again see fig 48 ('rullekake' eksempel.)
- But, if one attempts to represent data in a space with dimensionality lower than the intrinsic one, the problem of 'overcrowding' can occur. Means low-dimensional embedding of high-dimensional data are often ambiguous. = two points far apart in the data space are mapped to the vicinity of each other in the latent space.
- To alleviate this, one needs to weaken the constraint we impose on our visualization schemes. Fex: in the case of t-distributed stochastic embedding (t-SNE), one prioritizes the preservation of short distances or local ordination(?) rather than that of all pairwise distances.

12.2 Principal component analysis (PCA)

- Goal: perform a linear projection of the data onto a lower-dimensional subspace where the variance is maximized. Inspired by the observation that in many cases, relevant information is often contained in the directions with largest variance. (fig 50, helpful!)
- Intuitively, these dirs encode the large-scale 'signal' as opposed to 'noise' characterized by the dir of small variance.
- PCA also seeks variable dirs while simultaneously reducing the redundancy between new basis vectors. Done by requiring our new basis vectors (=principal components) be orthogonal.
- Data then visualized by projecting it onto a subspace spanned by a few principal component basis vectors.
- Surprisingly, such PCA-based projections often capture a lot of the large scale structure of many datasets. Fex, fig 51 (cool!): shows the projections of samples drawn from the 2D Ising model at various temperatures on the first two principal components. Despite living in a 1600 dimensional space (the samples are 40×40 spins), a single component (i.e. a single direction in this 1600 dimensional space) can capture 50% of the variability contained in our samples. One can actually check that easily that this direction weights all 1600 spins equally and thus **corresponds to the magnetization parameter**. → even without any prior physical knowledge, one can extract relevant order parameters using a simple PCA-based projection.

- PCA is widely employed in biological physics when working with high-dimensional data. Recently, **a correspondence between PCA and Renormalization Group flows across the phase transition in the 2D Ising model or in a general setting has been proposed.** In stat phys, PCA has also found application in detecting phase transitions, e.g. in the XY model on frustrated triangular and union jack lattices. Also used to classify dislocation patterns in crystals. Physics has also inspired PCA-based algos to infer relevant features in unlabelled data.
- Concretely,
 - Consider N data points $\{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ that live in a D -dimensional feature space \mathbb{R}^D .
 - Without loss of generality, we assume the empirical mean $\bar{\mathbf{x}} = N^{-1} \sum_i \mathbf{x}_i$ of these data points is zero (we can always center around the mean: $\bar{\mathbf{x}} = \mathbf{x}_i - \bar{\mathbf{x}}$).
 - Denote $N \times D$ design matrix as $X = [\bar{\mathbf{x}}_1, \bar{\mathbf{x}}_2, \dots, \bar{\mathbf{x}}_N]^T$ whose rows = the data points and columns = different features. The $D \times D$ (symmetric) covariance matrix is therefore

$$\Sigma(X) = \frac{1}{N-1} X^T X \quad (173)$$

Notice

- * The j -th diagonal entry of $\Sigma(X)$ = the variance of j -th feature
- * The $\Sigma(X)_{ij}$ measures the covariance (i.e. **connected correlation in the language of physics**) between feature i and j .
- We want: find a new basis for the data that emphasizes highly variable directions while reducing redundancy between basis vectors. In particular, we'll look for **a linear transformation that reduces the covariance between different features.** To do so,
 - * First, perform a singular value decomposition (SVD) on design matrix X , namely: $X = USV^T$, where S =a diagonal matrix of singular value s_i , the orthogonal matrix U contains (as its columns) the left singular vectors of X , and similarly V contains (as its columns) the right singular vectors of X .
 - * With this, we can rewrite the covariance matrix as

$$\Sigma(X) = \frac{1}{N-1} V S U^T U S V^T \quad (174)$$

$$= V \left(\frac{S^2}{N-1} \right) V^T = V \Lambda V^T \quad (175)$$

where Λ = a diagonal matrix w/ eigenvalues λ_i in the decreasing order along the diagonal (i.e. eigendecomposition).

- * Clear that the right singular vectors of X (=columns of V) are principle directions of $\Sigma(X)$, and
- * Singular values of X are related to the eigenvalues of covariance matrix $\Sigma(X)$ via $\lambda_i = s_i^2/(N - 1)$.
- * To reduce the dimensionality of data from D to $\tilde{D} < D$, first construct the $D \times \tilde{D}$ projection matrix $\tilde{V}_{D'}$ by selecting the singular components with the \tilde{D} largest singular values.
- * The projection of the data from D to a \tilde{D} dimensional space is simply $\tilde{Y} = X\tilde{V}_{D'}$.
- * The singular vector w/ the largest value (=largest variance) is referred to as the first principal component, the singular vector w/ the second largest singular value as the second principal component, etc.
- * An important quantity: the ratio $\lambda_i / \sum_{i=1}^D \lambda_i$ = the percentage of the explained variance contained in a principal component. (fig 51b)
- Common in data visualization to represent data projected on the first few principal components. Valid as long as a large part of the variance is explained in those components.
- Low values of explained variance may imply the intrinsic dimensionality of the data is high or simply that it cannot be captured by a linear representation.

12.3 Multidimensional scaling

- MDS = a non-linear dimensional reduction technique which preserves the pairwise distance/dissimilarity d_{ij} between data points. To types of MDS:
 - Metric MDS: the distance matrix is computed under a pre-defined metric and the latent coordinates \tilde{Y} are obtained by minimizing the difference between the distance matrix in the original space ($d_{ij}(X)$) and that in the latent space ($d_{ij}(Y)$):

$$\tilde{Y} = \operatorname{argmin}_Y \sum_{i < j} w_{ij} |d_{ij}(X) - d_{ij}(Y)| \quad (176)$$

where w_{ij} is a weight value: $w_{ij} \geq 0$. The weight matrix w_{ij} is a set of free parameters that specify the level of confidence (or precision) in the value of $d_{ij}(X)$. If Euclidean metric used, it's the same as PCA and usually called "classical scaling". Thus MDS often considered generalization of PCA.

- Non-metric MDS: d_{ij} can be any distance matrix. The objective function is then to preserve the ordination in the data. I.e. if $d_{12}(X) < d_{13}(X)$ in the original space, then in latent space we should have $d_{12}(X) < d_{13}(X)$.
- Both PCA and MDS can be implemented using standard Python packages such as Scikit. MDS algo typically typically have a scaling of $\mathcal{O}(N^3)$ where N =number of data points.
- PCA: if one is only interested in a small fraction of the principal components w/the largest variance (which is the usual case), efficient implementations based on Lanczos methods can achieve scaling of $\mathcal{O}(N^2)$ for dense matrices.
- PCA and MDS often among the first data visualization techniques one resorts to.

12.4 t-SNE

- Often desirable to preserve local structures in high-dimensional dataset. Typically not possible using linear techniques such as PCA.
- Many non-linear techniques such as non-classical MDS, self-organizing map, Isomap and Locally Linear Embedding have been proposed recently. These techniques generally good at preserving local structures in the data, but typically *fail to capture structures at the larger scale such as the clusters in which the data is organized*.
- Recently, t -stochastic neighbor embedding (t-SNE) has emerged as one of the go-to methods for visualizing high-dimensional data.
- t-SNE = a non-parametric method that utilizes non-linear embeddings. When used appropriately, a powerful technique for unraveling the hidden structures of high-dimensional datasets AND preserving locality.
- In **physics**, t-SNE has recently been used to reduce the dimensionality of and classify spin configs, generated with Monte-Carlo simulations, for the Ising and Fermi-Hubbard models at finite temps. Was also applied to study clustering transitions in random satisfiability problems which bears close resemblance to spin glass models.
- Idea of stochastic neighborhood embedding (SNE):
 - Associate a prob dist to the neighborhood of each data (note

$x \in \mathbb{R}^s$, s the number of features):

$$p_{i|j} = \frac{e^{-\frac{\|x_i - x_j\|^2}{2\sigma_i^2}}}{\sum_{k \neq i} e^{-\frac{\|x_i - x_k\|^2}{2\sigma_i^2}}} \quad (177)$$

where $p_{i|j}$ can be interpreted as the likelihood that x_j is x_i 's neighbor. σ_i are free parameters that are usually fixed by fixing the local entropy (=the perplexity) of each data point (regions of high density therefore have smaller σ_i).

- Intuitively, the Gaussian likelihood (i.e. short-tailed) means that only points that are nearby x_i contribute to its prob dist.
- A symmetrized prob dist is constructed from the above equation: $p_{ij} \equiv (p_{i|j} + p_{j|i})/(2N)$. The symmetrization ensures even outliers contribute p_{ij} and as such, have meaningful embedding coordinates.

- t -SNE, on the other hand

- Constructs an equivalent prob dist in a low dimensional latent space (with coordinates $y_i \in \mathbb{R}^t$, $t < s$, with t the dimension of the latent space):

$$q_{ij} = \frac{(1 + \|y_i - y_j\|^2)^{-1}}{\sum_{k \neq j} (1 + \|y_i - y_k\|^2)^{-1}} \quad (178)$$

The crucial point = q_{ij} is chosen to be a long-tail (Cauchy) distribution. This is meant to preserve short distance information (relative neighborhoods) while strongly repelling two points that are far apart in the original space (fig 52).

- In order to find the latent space coordinates y_i , t -SNE minimizes the Kullback-Leibler divergence between q_{ij} and p_{ij} :

$$KL(p||q) \equiv \sum_{ij} p_{ij} \log\left(\frac{p_{ij}}{q_{ij}}\right) \quad (179)$$

Achieve this by gradient descent.

- When visualizing data with t -SNE, important to understand what is plotted. Some important properties to bear in mind when analyzing t -SNE plots:
 - *t -SNE can rotate data* The KL divergence is invariant under rotations in the latent space, since it only depends on the distance between points. \rightarrow t -SNE plots that are rotations of each other should be considered equivalent.

- *t-SNE results are stochastic* Although KL divergence is convex in the domain of distributions, it is generally not in the domain of q_{ij} (i.e. latent coordinate y). \rightarrow in applying gradient descent the solution will depend on the initial seed. \rightarrow the map obtained may very depending on the seed used and different t-SNE runs will give slightly different results.
 - *t-SNE generally preserves short distance information* Rule of thumb: expect nearby points on the t-SNE points are also closeby in the original space. The reason is the nature of the mapping discussed in fig 52.
 - *Scales are deformed in t-SNE* Since a scale-free distribution is used in the latent space, one should not put too much emphasis on the meaning of the variances of any clusters observed in the latent space.
 - *t-SNE is computationally intensive* A direct implementation of t-SNE has an algo complexity of $\mathcal{O}(N^2)$ which is only applicable to small to medium data sets. Improved scaling of the form $\mathcal{O}(N \log N)$ can be achieved at the cost of approximating $KL(p||q)$ by using Barnes-Hut method.
- Illustration:
 - Fig 53: t-SNE applied to a model consisting of thirty Gaussians (=a Gaussian mixture model) whose means are uniformly distributed in a forty-dimensional space. Compared the results to a random 2D projection and PCA.
 - Clear that unlike more naive dimensional reduction techniques, both PCA and t-SNE can identify the presence of well-formed clusters. The t-SNE visualization cleanly separates all the clusters while certain clusters blend together in the PCA plot. = A direct consequence of the fact that t-SNE keeps nearby points close together while repelling points that are far apart.
 - Fig 54: t-SNE and PCA plots for MNIST dataset of ten hand-written numerical digits (0-9).
 - Clear that the non-linear nature of t-SNE makes it much better at capturing and visualizing the complicated correlations between digits than the PCA.

13 Clustering

- We continue discussion of unsupervised learning methods. Unsupervised learning concerned with: discovering structure in unlabeled data (fex learning local structures for data visualization).

- More difficult than supervised because unlabeled data. Surprising: still possible to uncover and exploit hidden structure in the data.
- Perhaps simplest example of unsupervised learning = clustering. Aim: group unlabeled data into clusters according to some similarity or distance measure. Informally thought of as: a set of points sharing some pattern or structure.
- Many applications of clustering in
 - Data mining
 - Data compression
 - Signal processing
 - Can be used to identify coarse features or high level structures in an unlabelled dataset.
 - Applications in physical sciences:
 - * Detecting celestial emission sources in astronomical surveys
 - * Inferring groups of genes and proteins with similar functions in biology
 - * Building entanglement classifiers
- Clustering=vast field, flurry of methods suited for different purposes. Common considerations when choosing method:
 - the distributions of the clusters (overlapping/noisy clusters vs well-separated clusters),
 - the geometry of the data (flat vs non-flat),
 - the cluster size distribution (multiple sizes vs uniform sizes),
 - dimensionality of the data (low vs high dimensional) and
 - the computational efficiency of the desired method (small vs large dataset).

13.1 Practical clustering methods

Throughout this section focus on Euclidean distance as similarity measure. See reference for more in depth discussion of different possible similarity measures.

13.1.1 K-means

- Consider a set of N *unlabeled* observations $\{\mathbf{x}_n\}_{n=1}^N$ where $\mathbf{x}_n \in \mathbb{R}^p$ and where p is the number of features.

- Also consider a set of K cluster centers called the cluster *means*: $\{\boldsymbol{\mu}_k\}_{k=1}^K$ with $\boldsymbol{\mu}_k \in \mathbb{R}^p$. The cluster means can be thought of as the representatives of each cluster, to which datapoints are assigned (fig 55).
- K -means clustering can be formulated as follows:

- Given a fixed integer K , find the cluster means $\{\boldsymbol{\mu}\}$ and the data point assignments in order to minimize the following objective function:

$$J(\{x, \boldsymbol{\mu}\}) = \sum_{k=1}^K \sum_{n=1}^N r_{nk} (\mathbf{x}_n - \boldsymbol{\mu}_k)^2 \quad (180)$$

where r_{nk} is a binary variable ($r_{nk} \in \{0, 1\}$) called the assignment.

- The assignment r_{nk} is 1 if x_n is assigned to cluster k and 0 otherwise. Notice that $\sum_k r_{nk} = 1 \quad \forall \quad n$ and $\sum_n r_{nk} \equiv N_k$, where N_k is the number of points assigned to cluster k .
- The minimization of this objective function can be understood as trying to find the best cluster means such that the variance within each cluster is minimized.
- In physical terms: J is the sum of the moments of inertia of every cluster. Indeed, as we'll see below, the cluster means $\boldsymbol{\mu}_k$ correspond to the centres of mass of their respective cluster.
- **K -means algorithm** K -means algo alternates between two steps:

1. *Expectation*: Given a set of assignments $\{r_{nk}\}$, minimize J wrt $\boldsymbol{\mu}_k$. Taking a simple derivative and setting it to zero yields the update rule:

$$\boldsymbol{\mu}_k = \frac{1}{N_k} \sum_n r_{nk} \mathbf{x}_n \quad (181)$$

2. *Maximization*: Given a set of cluster means $\{\boldsymbol{\mu}_k\}$, find the assignments $\{r_{nk}\}$ which minimizes J . Clearly, this is achieved by assigning each data point to their nearest cluster-mean:

$$r_{nk} = \begin{cases} 1, & \text{if } k = \operatorname{argmin}_{k'} (\mathbf{x}_n - \boldsymbol{\mu}_{k'})^2 \\ 0, & \text{otherwise} \end{cases}$$

K -means clustering consists in alternating between these two steps until some convergence criterion met. Practically it should terminate when the change in the objective func from one iteration to another becomes smaller than a pre-specified threshold. Fig 55: Simple example of algo.

- Nice property: K -means is guaranteed to converge. To see this: can verify explicitly (by taking second-order derivatives) that the expectation step always decreases J . Also true for the assignment step. \rightarrow since J bound from below, the two-step iteration *always* converges to a local min of J .
- J generally a non-convex func \rightarrow in practice one usually needs to run the algo with different initial random seeds and post-select the best local minimum.
- A simple implementation has an average computational complexity which scales linearly in the size of the data set (that is, $\mathcal{O}(KN)$ per iteration) and is thus scalable to very large datasets.
- As we'll see later, K -means is a hard-assignment limit of the Gaussian mixture model where all cluster variances are assumed to be the same.
- This highlights a common drawback of K -means: if the true clusters have very different variances (spreads), K -means can lead to spurious results since the underlying assumption is that the latent model has uniform variances.

13.1.2 Hierarchical clustering: Agglomerative methods

- Agglomerative clustering = A bottom-up approach that starts from small initial clusters which are then progressively merged to form larger clusters.
- The merging process generates a hierarchy of clusters that can be visualized in the form of a dendrogram (fig 56). This hierarchy can be useful to analyze the relation between clusters and the subcomponents of individual clusters.
- Agglomerative methods usually specified by defining a distance measure between clusters (the measure need not be metric). Denote the distance between clusters X and Y by $d(X, Y) \in \mathbb{R}$. Different choices of distance result in different clustering algos. At each step, the two clusters that are closest wrt the distance measure are merged until a single cluster is left.
- **Agglomerative clustering algo:**
 1. Initialize each point to its own cluster
 2. Given a set of K clusters X_1, X_2, \dots, X_K , merge clusters until one cluster is left ($K=1$):
 - (a) Find the closest pair of clusters $(X_i, X_j) : (i, j) = \operatorname{argmin}_{(i', j')} d(X_{i'}, X_{j'})$

(b) Merge the pair. Update $K \leftarrow K - 1$

- A few of the most popular distances used in agglomerative methods (often called linkage methods):

1. Single-linkage: the distance between clusters i and j = the minimum distance between two elements of the different clusters:

$$d(X_i, X_j) = \min_{\mathbf{x}_i \in X_i, \mathbf{x}_j \in X_j} \|\mathbf{x}_i - \mathbf{x}_j\|_2 \quad (182)$$

2. Complete linkage: the distance between clusters i and j = the maximum distance between two elements of the different clusters

$$d(X_i, X_j) = \max_{\mathbf{x}_i \in X_i, \mathbf{x}_j \in X_j} \|\mathbf{x}_i - \mathbf{x}_j\|_2 \quad (183)$$

3. Average linkage: average distance between points of different clusters

$$d(X_i, X_j) = \frac{1}{|X_i| \cdot |X_j|} \sum_{\mathbf{x}_i \in X_i, \mathbf{x}_j \in X_j} \|\mathbf{x}_i - \mathbf{x}_j\|_2 \quad (184)$$

4. Ward's linkage: Analogous to the K -means method as it seeks to minimize the total inertia. The distance measure is the "error squared" before and after merging which simplifies to:

$$d(X_i, X_j) = \frac{|X_i||X_j|}{|X_i \cup X_j|} (\boldsymbol{\mu}_i - \boldsymbol{\mu}_j) \quad (185)$$

- Common drawback of hierarchical methods: they don't scale well: at every step, a distance matrix between all clusters must be updated/computed. Efficient implementations achieve a typical computational complexity of $\mathcal{O}(N^2)$, making the method suitable for small to medium datasets.
- A simple but major speed-up: initialize the clusters with K -means using a large K (but still a small fraction of N), then proceed with hierarchical clustering. Advantage: this preserves the large-scale structure of the hierarchy while making use of the linear scaling of K -means. In this way hierarchical methods may be applied to very large datasets.

13.1.3 Density-based (DB) clustering

- Makes the intuitive assumption that clusters are defined by regions of space w/higher density of data points. Data points that constitute noise or that are outliers are expected to form regions of low density. This has the advantage of being able to consider clusters of multiple shapes and sizes while identifying outliers. Also suitable for large-scale applications.

- Core assumption: *relative* local density estimation of data is possible. = possible to order points according to their densities. Density estimates usually accurate for low-dimensional data but become unreliable for high-dimensional data due to large sampling noise.
- Here, we confine our discussion to one of the most widely used density clustering algos DBSCAN. We have also had great success with another recently introduced DB clustering variant which is similar in spirit. Also one of the authors has created a Python package which makes use of accurate density estimates via kernel methods combined with agglomerative clustering to produce fast and accurate density clustering (GitHub).
- **DBSCAN algo** DBSCAN = density-based spatial clustering of applications with noise.
 - Consider a set of N data points $X \equiv \{\mathbf{x}_n\}_{n=1}^N$.
 - We start by defining the ϵ -neighborhood of point \mathbf{x}_n as follows:

$$N_\epsilon(\mathbf{x}_n) = \{\mathbf{x} \in X | d(\mathbf{x}, \mathbf{x}_n) < \epsilon\} \quad (186)$$

$N_\epsilon(\mathbf{x}_n)$ are the data points at a distance smaller than ϵ from \mathbf{x}_n . As before, consider $d(\cdot, \cdot)$ to be the Euclidean metric (which yields spherical neighborhoods, fig 57) but other metrics may be better suited depending on the data.

- $N_\epsilon(\mathbf{x}_n)$ can be seen as a crude estimate of local density. \mathbf{x}_n is considered to be a *core-point* if a least **minPts** are in its ϵ -neighborhood. **minPts** = a free parameter of the algo that sets the scale of the size of the smallest cluster one should expect.
- Finally, a point \mathbf{x}_i = *density-reachable* if it is in the ϵ -neighborhood of a *core-point*.
- From these definitions, the algo can be simply formulated:
- \rightarrow Until all points in X have been visited; **do**
 - * Pick a point \mathbf{x}_i that has not been visited
 - * Mark \mathbf{x}_i as a visited point
 - * If \mathbf{x}_i is a core point; **then**
 - Find the set \mathcal{C} of all points that are *density-reachable* from \mathbf{x}_i .
 - \mathcal{C} now forms a cluster. Mark all points within that cluster as being visited.
- \rightarrow Return the cluster assignments $\mathcal{C}_1, \dots, \mathcal{C}_k$, with k =the number of clusters. Points that have not been assigned to a cluster are considered noise or out-liers.

- Note: DBSCAN doesn't require the user to specify the number of clusters but only ϵ and **minPts**. While it's common to heuristically fix these parameters, methods such as cross-validation can be used for their determination.
- Finally, note: DBSCAN very efficient since efficient implementations have a computational cost of $\mathcal{O}(N \log N)$.

13.2 Clustering and latent variables via the Gaussian mixture models

- We'll here approach clustering from a more abstract vantage point, and in the process, introduce many of the core ideas underlying supervised learning.
- Central concept in many unsupervised learning techniques: the idea of a latent or hidden variable. Though not directly observable, latent variables still influence the visible structure of the data. Fex: in clustering we can think of the cluster identity of each datapoint (=the cluster the data point belongs to) as a latent variable. Even though can't see the cluster label explicitly, we know points in same cluster tend to be closer together. The latent variables in our data (cluster identity) are a way of representing and abstracting the correlations between datapoints.
- In this language: can think of clustering as an algo to learn the most probable value of a latent variable (cluster identity) associated with each datapoint. Calc this latent variable requires additional assumptions about the structure of our data. Like all unsupervised learning algos, in clustering we must make an assumption about the underlying prob dist from which the data was generated.
- Our model for how the data is generated is often called our **generative model**. In clustering, we
 - assume that data points are assigned a cluster, with each cluster characterized by some cluster-specific prob dist (e.g. a Gaussian with some mean and variance that characterizes the cluster).
 - We then specify a procedure for finding the value of the latent variable. Often done by choosing the values of the latent variable that minimize some cost func.
- One common choice for a class of cost funcs for many unsupervised learning problems = Maximum Likelihood Estimation (MLE): choose the values of the latent variables that maximize the likelihood of the observed data under our generative model (= maximize the prob of

getting the observed data under our generative model). Such MLE equations often give rise to the kind of **Expectation Maximization (EM)** equations that we first encountered in the K -means clustering context.

- Gaussian Mixture Models (GMM) are a generative model often used in the context of clustering.

- In GMM, points are drawn from one of K Gaussians, each with its own mean $\boldsymbol{\mu}_k$ and covariance matrix $\boldsymbol{\Sigma}_k$,

$$\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) \sim e^{-\frac{(\mathbf{x}-\boldsymbol{\mu})\boldsymbol{\Sigma}^{-1}(\mathbf{x}-\boldsymbol{\mu})^T}{2}} \quad (187)$$

- Let's denote the prob a data point is drawn from mixture k by π_k . Then, the prob of generating a point \mathbf{x} in a GMM is given by

$$p(\mathbf{x}|\{\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k, \pi_k\}) = \sum_{k=1}^K \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \pi_k \quad (188)$$

- Given a dataset $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$, we can write the likelihood of the dataset as

$$p(\mathbf{X}|\{\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k, \pi_k\}) = \prod_{i=1}^N p(\mathbf{x}_i|\{\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k, \pi_k\}) \quad (189)$$

- For future reference, let's denote the set of parameters $\{\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k, \pi_k\}$ by $\boldsymbol{\theta}$.
- To see how we can use GMM and MLE to perform clustering, we introduce discrete binary K -dimensional latent variables \mathbf{z} for each data point \mathbf{x} whose k -th component is 1 if point \mathbf{x} was generated from the k -th Gaussian and zero otherwise (often called "one-hot variables"). Fex: if we're considering a Gaussian mixture with $K = 3$ we would have three possible values for $\mathbf{z} \equiv (z_1, z_2, z_3)$: $(1, 0, 0)$, $(0, 1, 0)$ and $(0, 0, 1)$.
- We cannot directly observe \mathbf{z} . It's a latent variable that encodes the cluster identity of point \mathbf{x} . Let's also denote all the N latent variables corresponding to dataset \mathbf{X} by \mathbf{Z} .
- Viewing the GMM as a generative model, we can write the prob $p(\mathbf{x}|\mathbf{z})$ of observing a data point \mathbf{x} given \mathbf{z} as

$$p(\mathbf{x}|\mathbf{z}; \{\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k\}) = \prod_{k=1}^K \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)^{z_k} \quad (190)$$

and the prob of observing a given value of latent variable

$$p(\mathbf{z}|\{\pi_k\}) = \prod_{k=1}^K \pi_k^{z_k} \quad (191)$$

- Using Bayes rule, we can write the joint prob of a clustering assignment \mathbf{z} and a data point \mathbf{x} given the GMM parameters as

$$p(\mathbf{x}, \mathbf{z}; \theta) = p(\mathbf{x}|\mathbf{z}; \{\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k\})p(\mathbf{z}|\{\pi_k\}) \quad (192)$$

- Can also use Bayes rule to rearrange this expression to give the conditional prob of the data point \mathbf{x} being in the k -th cluster, $\gamma(z_k)$, given model parameters θ as

$$\gamma(z_k) \equiv p(z_k = 1|\mathbf{x}; \theta) = \frac{\pi_k \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)} \quad (193)$$

$\gamma(z_k)$ are often referred to as the "responsibility" that mixture k takes for explaining \mathbf{x} . Just like in our discussion of soft-max classifiers, this can be made into a "hard-assignment" by assigning each point to the cluster with the largest prob: $\operatorname{argmax}_k \gamma(z_k)$ over the responsibilities.

- The complication is of course that we don't know the parameters θ of the underlying GMM but instead must also learn them from the data \mathbf{X} . As discussed, ideally we could do this by choosing the parameters that maximize the likelihood (or equivalently the log-likelihood) of the data

$$\hat{\theta}_i = \operatorname{argmax}_{\theta_i} \log p(\mathbf{X}|\theta) \quad (194)$$

Once we knew the MLEs $\hat{\theta}_i$, we can cal the optimal hard cluster assignment $\operatorname{argmax}_k \hat{\gamma}(z_k)$ where $\hat{\gamma}(z_k) = p(z_k = 1|\mathbf{x}; \hat{\theta})$.

- In practice, due to the complexity of the expression for $p(\mathbf{X}|\theta)$, it's almost impossible to find the global max of the likelihood func. Instead, must settle for local max. One approach to finding the local max is use a method like SGD on the negative log-likelihood.
- Here, we introduce an alternative, powerful approach for finding local minima in latent variable models using an iterative procedure called **Expectation Maximization (EM)**. Given an initial guess of the parameters $\theta^{(0)}$, the EM algo iteratively generates new estimates for the parameters $\theta^{(1)}, \theta^{(2)}, \dots$. Importantly, the likelihood is guaranteed to be non-decreasing under these iterations and hence EM converges to a local max of the likelihood.
 - Central observation underlying EM: often much easier to calc the conditional likelihoods of the latent variables $\tilde{p}(\mathbf{Z}) = p(\mathbf{Z}|\mathbf{X}; \theta^{(t)})$ given some choice of parameters and the max of the expected log likelihood given an assignment of the latent variables: $\theta^{(t+1)} = \operatorname{argmax}_{\theta} E_{p(\mathbf{Z}|\mathbf{X}; \theta^{(t)})} [\log p(\mathbf{X}, \mathbf{Z}; \theta)]$.

- To get an intuition for this latter quantity notice that we can write

$$E_{\tilde{p}^{(t)}}[\log p(\mathbf{X}, \mathbf{Z}; \theta)] = \sum_{i=1}^N \sum_{k=1}^K \gamma_{ik}^{(t)} [\log \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) + \log \pi_k] \quad (195)$$

where we used the shorthand $\gamma_{ik}^{(t)} = p(z_{ik} | \mathbf{X}; \theta^{(t)})$ with z_{ik} = the k -th component of \mathbf{z}_i

- Taking the derivative of this eq wrt $\boldsymbol{\mu}_k$, $\boldsymbol{\Sigma}_k$, and π_k (subject to the constraint $\sum_k \pi_k = 1$) and setting this to zero yields the intuitive equations

$$\boldsymbol{\mu}_k^{(t+1)} = \frac{\sum_i \gamma_{ik}^{(t)} \mathbf{x}_i}{\sum_i \gamma_{ik}^{(t)}} \quad (196)$$

$$\boldsymbol{\Sigma}_k^{t+1} = \frac{\sum_i \gamma_{ik}^{(t)} (\mathbf{x}_i - \boldsymbol{\mu}_k)(\mathbf{x}_i - \boldsymbol{\mu}_k)^T}{\sum_i \gamma_{ik}^{(t)}} \quad (197)$$

$$\pi_k^{(t+1)} = \frac{1}{N} \sum_k \gamma_{ik}^{(t)} \quad (198)$$

These are just the usual estimates for the mean and variance, with each data point weighed according to our current best guess for the prob that it belongs to cluster k .

- We can then use our new estimate $\theta^{(t+1)}$ to calc new memberships $\gamma_{ik}^{(t+1)}$ and repeat the process. This is essentially the K -Means algo discussed previously.
- This discussion of GMMs introduces several concepts that we'll return to repeatedly in the context of unsupervised learning.
 - First, it's often useful to think of the visible correlations between features in the data as resulting from hidden or latent variables.
 - Second, we'll often posit a generative model that encodes the structure we think exists in the data and then find parameters that maximize the likelihood of the observed data.
 - Third, often we'll not be able to directly estimate the MLE, and will have to instead look for a computationally efficient way to find a local min of the likelihood.

13.3 Clustering in high-dimension

- One major problem that is aggravated when clustering data in high-dimension = the accumulation of noise coming from spurious features

that tends to "blur" distances. Many clustering algos rely on the explicit use of a similarity measure or distance metrics that weigh all features equally. → one must be careful when using an off-the-shelf method in high dimensions.

- In order to perform clustering in high-dim data, often useful to denoise the data before proceeding w/using a standard clustering method.
 - Fig 54 (seen earlier): PCA was used to denoise the MNIST dataset by projecting the 748 original dimensions onto the 40 dimensions w/the largest principal components.
 - The resulting features were then used to construct a Euclidean distance matrix which was used by t -SNE to compute the 2D embedding that is presented. (Using t -SNE directly on original data leads to "blurring" of the clusters).
- However, simple feature selection or feature denoising (using *ex* PCA) can sometimes be insufficient for learning clusters due to the presence of large variations in the signal and noise of the features that are relevant for identifying the underlying clusters. (why is this a problem? aren't the features picked out by PCA the ones with large variance? or is there a different type of variance, which PCA somehow avoids (how, if so?)?) Recent promising work suggests one way to overcome these limitations is to learn the latent space *and* the cluster labels at the *same time*.
- We end the clustering section with a discussion on cluster validation—can be particularly hard for high-dim data. Often cluster validation (= verifying whether the obtained labels are "valid") is done by direct visual inspection. That is: the data is represented in a low-dim space and the cluster labels obtained are visually inspected to make sure different labels organize into distinct "blobs".
- For high-dim data, this is done by performing dim-reduction. But, this can lead to the appearance of spurious clusters since dim-reduction inevitably loses info about the original data. → these methods should be used with care when trying to validate clusters (see ref for an interactive discussion on how t -SNE can sometime be misleading and how to effectively use it).
- There's a lot of work done devising ways of validating clusters based on various metrics and measures. Perhaps one of the most intuitive ways of defining a good clustering = measuring how well clusters generalize. Recently, clustering methods based on leveraging powerful classifiers to measure the generalization errors of the clusters have

been developed by some of the authors and we believe this represent an especially promising research direction for high-dim clustering.

- Finally, see ref for an in-depth survey of the various validation techniques.

14 Variational methods and mean-field theory (MFT)

- Common thread in many supervised learning tasks = accurately representing the underlying prob dist from which a dataset is drawn. Unsupervised learning of high-dim, complex distributions represents a new set of technical and computational challenges that are different from those we encountered in supervised learning.
- When dealing with complicated prob dists, often much easier to learn the *relative weights* of different states or data points (ratio of probs), than *absolute* probs.
- In **physics**, this is the familiar statement that the weights of a Boltzmann dist are much easier to calc than the partition func. The relative prob of two configs \mathbf{x}_1 and \mathbf{x}_2 are prop to the difference between their Boltzmann weights

$$\frac{p(\mathbf{x}_1)}{p(\mathbf{x}_2)} = e^{-\beta(E(\mathbf{x}_1) - E(\mathbf{x}_2))} \quad (199)$$

where as is usual in stat mech $\beta =$ the inverse temp and $E(\mathbf{x}; \theta) =$ the energy of state \mathbf{x} given some parameters (couplings) θ .

- However, calc the absolute weight of a config requires knowledge of the partition func

$$Z_p = \text{Tr}_{\mathbf{x}} e^{-\beta E(\mathbf{x})} \quad (200)$$

(where the trace is taken over all possible configs \mathbf{x}) since

$$p(\mathbf{x}) = \frac{e^{-\beta E(\mathbf{x})}}{Z_p} \quad (201)$$

In general, calc the partition func $Z_p =$ analytically and computationally intractable.

- Fex: for the Ising model with N binary spins, the trace involves calc a sum over 2^N terms, whcih is a difficult task for most energy funcs. \rightarrow physicists (and ML scientists) have developed various numerical and computational methods for evaluating such partition funcs.

- One approach: use Monte-Carlo based methods to draw samples from the underlying dist (can be known knowing only the relative probs) and then use these samples to numerically estimate the partition func.
- This is the philosophy behind powerful methods inc Markov Chain Monte Carlo (MCMC) and annealed importance sampling which are widely used in both the stat phys and ML communities.
- Alternative approach, that we focus on here: approximate the prob dist $p(\mathbf{x})$ and partition func using a "variational distribution" $q(\mathbf{x}; \theta_q)$ whose partition func we can calc exactly. The variational parameters θ_q are chosen to make the variational dist as close to the true dist as possible (how this is done = focus of much of this chapter).
- One of the most widely applied examples of a variational method in stat phys = Mean-Field Theory (MFT). MFT can be naturally understood as a procedure for approximating the true dist of the system by a **factorized distribution**. The deep connection between MFT and variational methods is discussed below. These variational MFT methods have been extended to understand more complicated spin models (aka graphical models in the ML literature) and form the basis of powerful set of techniques called **Belief Propagation** and **Survey Propagation**.
- Variational methods also widely used in ML to approximate complex probabilistic models. Fex: below we show how the EM procedure, discussed previously in the GMM clustering context, is actually a general method that can be dervied for any latent (hidden) variable model using a variational procedure.
- For readers interested in in-depth discussion on **variational inference for probabilistic graphical models**, we recommend a great treatise + a physics oriented discussion + an outstanding book **REFERENCES**

14.1 Variational mean-field theory for the Ising model

- Ising models = a major paradigm in stat phys
- Historically introduced to study magnetism, it was quickly realized their predictive power applies to a variety of interacting many-particle systems. Ising models are now understood to serve as minimal models for complex phenomena such as certain classes of phase transitions.
- In the model, degrees of freedom called spins assume discrete, binary values, e.g. $s_i = \pm 1$. Each spin variable s_i lives on a lattice (or in general, a graph), the sites of which are labeled by $i = 1, 2, \dots, N$.

- Despite the extreme simplicity relative to real-world systems, Ising models exhibit a high level of intrinsic complexity, and the degrees of freedom can become correlated in sophisticated ways. Often, spins interact spatially locally, and respond to externally applied magnetic fields.
- A spin config \mathbf{s} specifies the values s_i of the spins at every lattice site. We can assign an "energy" to every such config

$$E(\mathbf{s}, \mathbf{J}) = -\frac{1}{2} \sum_{i,j} J_{ij} s_i s_j - \sum_i h_i s_i \quad (202)$$

where h_i =a local magnetic field acting on the spin s_i , and J_{ij} = the interaction strength between the spin s_i and s_j . In textbook examples, the coupling parameters $\mathbf{J} = (J, h)$ =typically uniform or, in studies of disordered systems, drawn from some prob dist (i.e. quenched disorder).

- The prob of finding the system in a given spin config at temp β^{-1} is given by

$$p(\mathbf{s}|\mathbf{J}) = \frac{1}{Z_p(\mathbf{J})} e^{-\beta E(\mathbf{s}, \mathbf{J})} \quad (203)$$

$$Z_p(\mathbf{J}) = \sum_{s_i=\pm 1} e^{-\beta E(\mathbf{s}, \mathbf{J})} \quad (204)$$

with $\sum_{s_i=\pm 1}$ denoting the sum over all possible configs of the spin variables. We write Z_p to emphasize that it's the partition func corresponding to the prob dist $p(\mathbf{s})$ (will become important later).

- For a fixed number of lattice sites N , there are 2^N possible configs, a number that grows exponentially w/the system size. \rightarrow not in general feasible to evaluate the partition func $Z_p(\mathbf{J})$ in closed form. = a major obstacle for extracting predictions from physical theories since the partition func is directly related to the free-energy through the expression

$$\beta F_p(\mathbf{J}) = -\log Z_p(\mathbf{J}) = \beta \langle E(\mathbf{s}, \mathbf{J}) \rangle_p - H_p \quad (205)$$

with

$$H_p = - \sum_{s_i=\pm 1} p(\mathbf{s}|\mathbf{J}) \log p(\mathbf{s}|\mathbf{J}) \quad (206)$$

the entropy of the prob dist $p(\mathbf{s}|\mathbf{J})$.

- Even if the true prob $p(\mathbf{s}|\beta, \mathbf{J})$ may be a very complicated object, we can still make progress by approximating it by a *variational distribution* $q(\mathbf{s}, \boldsymbol{\theta})$ which captures the essential features of interest, with $\boldsymbol{\theta}$ some parameters that define our variational ansatz.
- The name variational dist comes from the fact that we're going to vary the parameters $\boldsymbol{\theta}$ to make $q(\mathbf{s}, \boldsymbol{\theta})$ as close to $p(\mathbf{s}|\beta, \mathbf{J})$ as possible.
- The functional form of $q(\mathbf{s}, \boldsymbol{\theta})$ = based on an "educated guess", often-times coming from our intuition about the problem.
- Can also define the variational free-energy

$$\beta F_q(\mathbf{J}, \boldsymbol{\theta}) = \beta \langle E(\mathbf{s}, \mathbf{J}) \rangle_q - H_q \quad (207)$$

where $\langle E(\mathbf{s}, \mathbf{J}) \rangle_q$ = the expectation value of the energy corresponding to the dist $p(\mathbf{s})$ wrt the dist $q(\mathbf{s}, \boldsymbol{\theta})$, and H_q is the entropy of $q(\mathbf{s}, \boldsymbol{\theta})$.

- Before proceeding further, helpful to introduce a new quantity: the Kullback-Leibler divergence (=KL-divergence=relative entropy) between two dists $p(\mathbf{x})$ and $q(\mathbf{x})$. It measures the dissimilarity between the two dists and is given by

$$D_{KL}(q||p) = \text{Tr} q(\mathbf{x}) \log \frac{p(\mathbf{x})}{q(\mathbf{x})} \quad (208)$$

= the expectation wrt q of the logarithmic difference between the two dists p and q . Two important properties of the KL-divergence are:

1. Positivity: $D_{KL}(p||q) \geq 0$ with equality iff $p = q$ (in the sense of prob dists)
 2. $D_{KL}(p||q) \neq D_{KL}(q||p)$, that is the KL-divergence is not symmetric in its arguments.
- Variational mean-field theory is a systematic way of constructing such an approximate distribution $q(\mathbf{s}, \boldsymbol{\theta})$. Main idea: choose parameters that minimize the difference between the variational free-energy $F_q(\mathbf{J}, \boldsymbol{\theta})$ and the true free energy $F_q(\mathbf{J}|\beta)$. Will show in following section that the difference between these two free-energies is actually the KL-divergence:

$$F_q(\mathbf{J}, \boldsymbol{\theta}) = F_q(\mathbf{J}|\beta) + D_{KL}(q||p) \quad (209)$$

This, combined w/the non-negativity of the KL-divergence has important consequences.

- First, it shows the variational free-energy is always larger than the true free-energy, $F_q(\mathbf{J}, \boldsymbol{\theta}) \geq F_p(\mathbf{J})$, with equality only if $q = p$ (the latter inequality is found in many physics textbooks and is known as the **Gibbs inequality**).
- Second, finding the best variational free-energy is equivalent to minimizing the KL divergence $D_{KL}(q||p)$.
- **Armed with these observations, let's derive a MFT of the Ising model using variational methods.** In the simplest MFT of the Ising model, the variational dist is chosen so that all spins are independent:

$$q(\mathbf{s}, \boldsymbol{\theta}) = \frac{1}{Z_q} e^{\sum_i \theta_i s_i} = \prod_i \frac{e^{\theta_i s_i}}{2 \cosh \theta_i} \quad (210)$$

Aka we've chosen a dist q which factorizes on every lattice site. An important property of this functional form = we can analytically find a closed-form expression for the variational partition func Z_q . This simplicity also comes at a cost= ignoring correlations between spins. These correlations become less and less important in higher dims, and the MFT ansatz becomes more accurate.

- To evaluate the variational free-energy, we
 - First need the entropy H_q . Since q factorizes over the lattice sites, the entropy separates into a sum of one-body terms

$$H_q(\boldsymbol{\theta}) = - \sum_{s_i = \pm 1} q(\mathbf{s}, \boldsymbol{\theta}) \log q(\mathbf{s}, \boldsymbol{\theta}) \quad (211)$$

$$= - \sum_i q_i \log q_i + (1 - q_i) \log(1 - q_i) \quad (212)$$

where $q_i = \frac{e^{\theta_i}}{2 \cosh \theta_i}$ is the prob that spin s_i is in the +1 state.

- Next, need to evaluate the average of the Ising energy $E(\mathbf{s}, \mathbf{J})$ wrt the variational dist q . Although the energy contains bilinear terms, we can still evaluate this average easily, because the spins are independent (uncorrelated) in the q dist. The mean value of spin s_i in the q dist, or on the on-site magnetization, is given by

$$m_i = \langle s_i \rangle_q = \sum_{s_i = \pm 1} s_i \frac{e^{\theta_i s_i}}{2 \cosh \theta_i} = \tanh(\theta_i) \quad (213)$$

Since the spins are independent, we have

$$\langle E(\mathbf{s}, \mathbf{J}) \rangle_q = -\frac{1}{2} \sum_{i,j} J_{ij} m_i m_j - \sum_i h_i m_i \quad (214)$$

The total variational free-energy is

$$\beta F_q(\mathbf{J}, \boldsymbol{\theta}) = \beta \langle E(\mathbf{s}, \mathbf{J}) \rangle_q - H_q \quad (215)$$

and minimizing wrt the variational parameters $\boldsymbol{\theta}$, we obtain

$$\frac{\partial}{\partial \theta_i} \beta F_q(\mathbf{J}, \boldsymbol{\theta}) = 2 \frac{dq_i}{d\theta_i} (-\beta [\sum_j J_{ij} m_j + h_i] + \theta_i) \quad (216)$$

Setting this eq to zero, we arrive at

$$\theta_i = \beta \sum_j J_{ij} m_j(\theta_j) + h_i \quad (217)$$

- For the special case of a uniform field $h_i = h$ and uniform nearest neighbor couplings $J_{ij} = J$, by symmetry the variational parameters for all the spins are identical, with $\theta_i = \theta$ for all i . Then, the mean-field equations reduce to their familiar textbook form, $m = \tanh(\theta)$ and $\theta = \beta(zJm(\theta) + h)$, where z is the coordination number of the lattice (i.e. the number of nearest neighbors).
- The equations for m_i and θ_i form a closed system = the mean-field equations for the Ising model. To solve them, one method is to iterate through and update each θ_i , once at a time, in an asynchronous fashion. To see the relationship of this approach to solving MFT equations to EM, it's helpful to explicitly spell out the iterative procedure to find the solutions to the eq for θ_i .
 - We start by initializing our variational parameters to some $\boldsymbol{\theta}^{(0)}$ and repeat the following until convergence:
 1. *Expectation*: Given a set of assignments at iteration t , $\boldsymbol{\theta}^{(t)}$, calc the corresponding magnetizations $\mathbf{m}^{(t)}$.
 2. *Maximization*: Given a set of magnetizations m_t , find new assignments $\theta^{(t+1)}$ which maximize the variational free-energy F_q . From our found expression for θ_i , this is just

$$\theta_i^{(t+1)} = \beta \sum_j J_{ij} m_j^{(t)} + h_i \quad (218)$$

From these equations, clear that we can think of the MFT of the Ising model as an EM like procedure similar to the one we used for k-means clustering and GMMs.

- As is well known in physics, even though MFT not exact, it can often yield qualitatively and even quantitatively precise predictions (especially in high dims). The discrepancy between the true physics and

MFT predictions stems from the fact that the variational dist q we chose doesn't model the correlations between spins. Fex: it predicts the wrong value for the critical temp for the 2D Ising model. Even erroneously predicts the existence of a phase transition in 1D at a non-zero temp.

- But we emphasize that the failure of any particular ansatz doesn't compromise the power of the approach. In some cases, one can consider changing the variational ansatz to improve the predictive properties of the corresponding variational MFT.

14.2 Expectation Maximization (EM)

- Ideas along the lines of variational MFT have been independently developed in statistics and imported into ML to perform maximum likelihood estimation (MLE).
- Here, we explicitly derive the EM algo and demonstrate further its close relation to MFT (**article** by Hinton).
- We'll focus on latent variable models where some of the variables=hidden=cannot be directly observed. Often makes MLE difficult to implement. EM gets around this difficulty by using an iterative two-step procedure, closely related to variational free-energy based approximation schemes in stat phys.
- Let \mathbf{x} = the set of visible variables we can directly observe and \mathbf{z} = the set of latent/hidden variables we cannot directly observe.
- Denote underlying prob dist from which \mathbf{x} and \mathbf{z} drawn by $p(\mathbf{z}, \mathbf{x}|\boldsymbol{\theta})$, with $\boldsymbol{\theta}$ = all relevant parameters.
- Given a dataset \mathbf{x} , we want to find the MLE of the parameters $\boldsymbol{\theta}$ that maximizes the prob of the observed data.
- As in variational MFT, we view $\boldsymbol{\theta}$ as variational parameters chosen to maximize the log-likelihood $L(\boldsymbol{\theta}) = \log p(\mathbf{x}|\boldsymbol{\theta})$. Algorithmically, this can be done by iterating the variational parameters $\boldsymbol{\theta}^{(t)}$ in a series of steps ($t = 1, 2, \dots$) starting from some arbitrary initial value $\boldsymbol{\theta}^{(0)}$:

1. **Expectation step (E step):** Given the known values of observed variable \mathbf{x} and the current estimate of parameter $\boldsymbol{\theta}_{t-1}$, find the prob dist of the latent variable \mathbf{z} :

$$q_{t-1}(\mathbf{z}) = p(\mathbf{z}|\boldsymbol{\theta}^{(t-1)}, \mathbf{x}) \quad (219)$$

2. **Maximation step (M step):** Re-estimate the parameter $\theta^{(t)}$ to be those with max likelihood, assuming $q_{t-1}(\mathbf{z})$ found in the prev step is the true dist of hidden variable \mathbf{z} :

$$\theta_t = \operatorname{argmax}_{\theta} \langle \log p(\mathbf{z}, \mathbf{x} | \theta) \rangle_{q_{t-1}} \quad (220)$$

It's been shown that each EM iteration increases the true log-likelihood $L(\theta)$, or at worst leaves it unchanged. In most models, this iteration procedure converges to a *local maximum* of $L(\theta)$.

- To see how EM is actually performed and related to variational MFT, we make use of the KL-divergence. Recall our goal =to maximize the log-likelihood $L(\theta) = \log p(\mathbf{x} | \theta)$. With data \mathbf{z} missing, we surely cannot just maximize $L(\theta)$ directly since parameter θ might couple both \mathbf{z} and \mathbf{x} . EM circumvents this by optimizing another objective func, $F_q(\theta)$, constructed based on estimates of the hidden variable dist $q(\mathbf{z} | \mathbf{x})$. Indeed, the func optimized is none other than the *variational free energy* encountered in the prev section:

$$F_q(\theta) := -\langle \log p(\mathbf{z}, \mathbf{x} | \theta) \rangle_q - H_q \quad (221)$$

where H_q is the SHannon entropy of $q(\mathbf{z})$. Can define the true free-energy $F_p(\theta)$ as the negative log-likelihood of the observed data:

$$-F_p(\theta) = L(\theta) = \log p(\mathbf{x} | \theta) \quad (222)$$

In the language of stat phys, $F_p(\theta)$ =the *true* free-energy while $F_q(\theta)$ =the variational free-energy we would like to minimize. Note we have employed a physics sign convention here of defining the free-energy as minus log of the partition func. In the ML literature, this minus often omitted, can lead to some confusion.

- Our goal: choose θ so that our variational free-energy $F_q(\theta)$ =as close to the true free-energy $F_p(\theta)$ as possible. The difference of these free-energies can be written as

$$F_q(\theta) - F_p(\theta) \quad (223)$$

$$= \log p(\mathbf{x} | \theta) - \sum_{\mathbf{z}} q(\mathbf{z} | \mathbf{x}) \log p(\mathbf{z}, \mathbf{x} | \theta) + \sum_{\mathbf{z}} q(\mathbf{z} | \mathbf{x}) \log q(\mathbf{z} | \mathbf{x}) \quad (224)$$

$$= \sum_{\mathbf{z}} q(\mathbf{z} | \mathbf{x}) \log p(\mathbf{x} | \theta) - \sum_{\mathbf{z}} q(\mathbf{z} | \mathbf{x}) \log p(\mathbf{z}, \mathbf{x} | \theta) + \sum_{\mathbf{z}} q(\mathbf{z} | \mathbf{x}) \log q(\mathbf{z} | \mathbf{x}) \quad (225)$$

$$= - \sum_{\mathbf{z}} q(\mathbf{z} | \mathbf{x}) \log \frac{p(\mathbf{z}, \mathbf{x} | \theta)}{p(\mathbf{x} | \theta)} + \sum_{\mathbf{z}} q(\mathbf{z} | \mathbf{x}) \log \tilde{p}(\mathbf{z}) \quad (226)$$

$$= \sum_{\mathbf{z}} q(\mathbf{z} | \mathbf{x}) \log \frac{q(\mathbf{z} | \mathbf{x})}{p(\mathbf{z} | \mathbf{x}, \theta)} \quad (227)$$

$$= D_{KL}(q(\mathbf{z} | \mathbf{x}) || p(\mathbf{z} | \mathbf{x}, \theta)) \geq 0 \quad (228)$$

where we've used Baye's theorem $p(\mathbf{z}|\mathbf{x}, \boldsymbol{\theta}) = p(\mathbf{z}, \mathbf{x}|\boldsymbol{\theta})/p(\mathbf{x}|\boldsymbol{\theta})$. Since KL-divergence always positive, this shows the variational free-energy F_q always an upper bound of the true free-energy F_p . In physics: **Gibb's inequality**.

- From eq for $F_q(\boldsymbol{\theta})$ and the fact that the entropy term there doesn't depend on **theta**, we can immediately see that the M-step=minimizing the variational free-energy $F_q(\boldsymbol{\theta})$.
- Surprisingly, the E-step can also be viewed as the optimization of the variational free-energy. Concretely, one can show that the dist of hidden variables \mathbf{z} given the observed variable \mathbf{x} and the current estimate of parameter $\boldsymbol{\theta}$ is the *unique* prob $q(\mathbf{z})$ that minimizes $F_q(\boldsymbol{\theta})$ (now seen as a functional of q). This can be proved by: taking the functional derivative of the equation for $F_q(\boldsymbol{\theta})$, plus a Langrange multiplier that encodes $\sum_{\mathbf{z}} g(\mathbf{z}) = 1$, wrt $q(\mathbf{z})$.
- Summing things up, we can rewrite EM as (same ref to Hinton):
 1. *Expectation step*: Construct the approximating prob dist of unobserved \mathbf{z} given the values of observed variable \mathbf{z} given the values of observed variable \mathbf{x} and parameter estimate $\boldsymbol{\theta}^{(t-1)}$:

$$q_{t-1}(\mathbf{z}) = \operatorname{argmin}_q F_q(\boldsymbol{\theta}^{(t-1)}) \quad (229)$$

2. *Maximation step*: Fix q , update the variational parameters:

$$\boldsymbol{\theta}^{(t)} = \operatorname{argmax}_{\boldsymbol{\theta}} -F_{q_{t-1}}(\boldsymbol{\theta}) \quad (230)$$

- **To recapitulate:**
 - EM implements MLE even with missing or hidden variables through optimizing a lower bound of the true log-likelihood.
 - In statistical physics, this reminiscent of optimizing a variational free-energy which is a lower bound of true free-energy due to Gibbs inequality.
- Fig 59: picture of how EM works
 - The E-step can be seen as representing the unobserved variable \mathbf{z} by a prob dist $q(\mathbf{z})$. This prob used to construct an alternative objective func $-F_q(\boldsymbol{\theta})$, which is then maximized wrt $\boldsymbol{\theta}$ in the M-step.
- By construction, maximizing the negative variational free-energy is equivalent to doing MLE on the joint data (=both observed and unobserved).

- "M-step" name = intuitive since θ found by maximizing $-F_q(\theta)$. "E-step" comes from the fact that one usually doesn't need to construct the prob of missing datas explicitly, but rather need only to compute the "expected" sufficient statistics over these data.
- Practical side: Em demonstrated as extremely useful in parameter estimation (but we're doing prediction in ML...?), particularly in hidden Markov models and Bayesian networks. Striking advantage: conceptually simple + easy to implement. In many cases, implementation of EM guaranteed to increase the likelihood monotonically = could be a perk during debugging.
- For readers familiar with physics of disordered systems: it's possible to construct a one-to-one dictionary between EM for latent variable models and the MFT of spin systems with quenched disorder. In disordered spin systems, the Ising couplings \mathbf{J} commonly taken to be quenched random variables drawn from some underlying prob dist. In the EM procedure, the quenched disorder provided by the observed data \mathbf{x} which are drawn from some underlying prob dist characterizing the data. The spins \mathbf{s} are like the hidden or latent variables \mathbf{z} . Similar analogies can be found for all the variational MFT quantities (table 1):

- **Statistical physics - Variational EM**
- spins/d.o.f: $\mathbf{s} \Leftrightarrow$ Hidden/latent variables \mathbf{z}
- couplings/quenched disorder: $\mathbf{J} \Leftrightarrow$ Data observations: \mathbf{x}
- Boltzmann factor $e^{-\beta E(\mathbf{s}, \mathbf{J})} \Leftrightarrow$ Complete probability: $p(\mathbf{x}, \mathbf{z} | \theta)$
- Partition function: $Z(\beta, \mathbf{J}) \Leftrightarrow$ Marginal likelihood $P(\mathbf{x} | \theta)$
- Energy: $\beta E(\mathbf{s}, \mathbf{J}) \Leftrightarrow$ Negative log-complete data likelihood: $-\log p(\mathbf{x}, \mathbf{z} | \theta, m)$
- Free energy: $\beta F_p(\mathbf{J} | \beta) \Leftrightarrow$ negative log-marginal likelihood: $-\log p(\mathbf{x} | m)$
- Variational distribution: $q(\mathbf{s}) \Leftrightarrow$ Variational distribution: $q(\mathbf{z} | \mathbf{x})$
- Variational free-energy: $F_q(\mathbf{J}, \theta) \Leftrightarrow$ Variational free-energy: $F_q(\mathbf{x}, \theta)$

15 Energy based models: Maximum Entropy (Max-Ent) Principle, Generative models, and Boltzmann Learning

- *Discriminative* models:
 - Most of models we have discussed so far, inc linear and logistic regression, ensemble models, supervised neural networks
 - Designed to perceive differences between groups/categories of data.

- Fex: recognizing differences between images of cats and images of dogs allows a discriminative model to label an image as "cat" or "dog".
- Form the core techniques of most supervised learning methods.
- Several limitations:
 1. Like all supervised learning, require labeled data.
 2. There are tasks that discriminative approaches simply cannot accomplish, fex: drawing new examples from an unknown prob dist.
- *Generative* models:
 - A model that can learn to represent and sample from a prob dist
 - Fex: A gen model for images would learn to draw new examples of cats and dogs given a dataset of images of cat and dog. And: Given samples generated from one phase of an Ising model we may want to generate new samples from that phase.
 - Such tasks clearly beyond the scope of discriminative models like the ensemble methods and DNNs discussed so far.
- *Energy-based* generative models:
 - Closely related to the kinds of models commonly encountered in statistical physics → we'll draw on many techniques that have their origin in stat mech (e.g. Monte-Carlo methods)
- Chapter overview:
 - Overview of generative models, highlighting similarities and differences with the supervised learning methods previously encountered
 - Introduce perhaps simplest kind of generative model: Maximum Entropy (MaxEnt) models. Have no latent/hidden variables=ideal for introducing key concepts underlying energy-based generative models.
 - Extended discussion on how to train energy-based models. Much of it also applicable to more complicated energy-based models such as the RBM and the deep models discussed later.

15.1 An overview of energy-based generative models

- Generative models = ML methods that learn to generate new examples similar to those found in a training dataset. Core idea for most: learn a parametric model for the prob dist from which the data was drawn. Can then generate new examples by sampling from the learned model.

- In stat phys, this sampling often done using Markov Chain Monte Carlo (MCMC) methods. **A concise and beautiful intro to MCMC-inspired methods that bridges both stat phys and ML is REFs (book and a review).**
- Added complexity of learning models directly from samples introduces many of the same fundamental tensions encountered when discussing discriminative models:
 - Model must be able to "generalize" beyond the examples they have been trained on = generate new samples not in the training set
 - → must be expressive enough to capture complex correlations present in the underlying data distribution
 - → amount of data we have is finite, giving rise to overfitting
- In practice: most generative models used in ML flexible enough that, with sufficient number of parameters, they can approximate any prob dist → there are thus **three axes on which we can differentiate classes of generative models**:
 1. How easy the model is to train - in terms of
 - Computational time
 - The complexity of writing code for the algo
 2. How well the model generalizes from the training set to the test set
 3. **Which characteristics of the data dist the model is capable of and focuses on capturing**
- One of fundamental reasons energy-based models less widely-employed than their discriminative counterpart is the training procedures of these models differ significantly from those for supervised NN models:
 - Both employ gradient-descent based procedures for minimizing a cost function (one common choice for generative models=the negative log-likelihood func)
 - But: energy-based models do not use backprop and automatic differentiation for computing gradients
 - They use: ideas inspired by MCMC based methods in physics and statistics, sometimes named "Boltzmann learning"
 - → Training them requires additional tools not immediately available in packages like PyTorch and TensorFlow

- Paysage, built on top of PyTorch bridges this gap, provides training for methods like RBM and stacked RBMS. Can be employed on GPUs. Maintained by a company affiliated with the authors.
- Finally: generative models at their most basic level = complex parametrizations of the prob dist the data is drawn from. → can do much more than just generate new examples. They can be used to perform a multitude of other tasks that require sampling from a complex prob dist including:
 - "de-noising"
 - filling in missing data
 - and even discrimination (Hinton reference).

The versatility is one of the major appeals of generative models.

15.2 Maximum entropy models: the simplest energy-based generative models

- Have their origin in a series of beautiful papers by Jaynes that reformulated stat mech in information theoretic terms **REF**.
- Often presented as the class of generative models that make the least assumptions about the underlying data - BUT all models make assumptions.

15.2.1 MaxEnt models in statistical mechanics

- Introduced by E. T. Jaynes in a two-part paper in 1957 entitled "**Information theory and statistical mechanics**". In these incredible papers, Jaynes
 - Showed it was possible to rederive the Boltzmann dist (and the idea of generalized ensembles) entirely from information theoretic arguments.
 - Quoting from the abstract, Jaynes considered "**stat mech as a form of stat inference rather than as a physical theory**" (portending the close connection between statistical physics and ML).
 - Showed the Boltzmann dist could be viewed as resulting from a stat inference procedure for learning prob dists describing physical systems where one only has partial information about the system (usually the average energy).

- The key quantity in MaxEnt models = the information theoretic entropy = Shannon entropy = a concept introduced by Shannon in his landmark treatise on information theory (1949).
 - It quantifies the stat uncertainty one has about the value of a random variable \mathbf{x} drawn from a prob dist $p(\mathbf{x})$.
 - It's defined as

$$S_p = -\text{Tr}_{\mathbf{x}} p(\mathbf{x}) \log p(\mathbf{x}) \quad (231)$$

where the trace = a sum/integral over all possible values a variable can take.

- Janyes showed the Boltzmann dist follows from the Principle of Maximum Entropy.
 - A physical system should be described by the prob dist w/the largest entropy subject to certain constraints (often provided by measuring the average value of conserved, extensive quantities such as the energy, particle number, etc.)
 - The principle uniquely specifies a procedure for parameterizing the functional form of the prob dist.
 - Having specified and learned this form we can, ofc, generate new examples by sampling this dist.
- More details on the workings:
 - $\{f_i(\mathbf{x})\}$ = set of funcs we have chosen whose average value we want to fix to some observed values $\langle f_i \rangle_{obs}$.
 - Principle of Max Ent states we should choose the dist $p(\mathbf{x})$ w/the largest uncertainty (i.e. largest Shannon entropy S_p), subject to the constraints that the model averages match the observed averages:

$$\langle f_i \rangle_{model} := \int d\mathbf{x} f_i(\mathbf{x}) p(\mathbf{x}) = \langle f_i \rangle_{obs} \quad (232)$$

- Can formulate the Principle of Max Ent as an optimization problem using the method of Lagrange multipliers (=the λ_i below) by minimizing:

$$\mathcal{L}[p] = -S_p + \sum_i \lambda_i (\langle f_i \rangle_{obs} - \int d\mathbf{x} f_i(\mathbf{x}) p(\mathbf{x})) + \gamma (1 - \int d\mathbf{x} p(\mathbf{x})) \quad (233)$$

where the first set of constraints=enforce the requirement for the averages and the last constraint = enforces the normalization that the trace over the prob dist equals one.

- Can solve for $p(\mathbf{x})$ by taking the functional derivative and setting it to zero

$$0 = \frac{\delta \mathcal{L}}{\delta p} = (\log p(\mathbf{x}) + 1) - \sum_i \gamma_i f_i(\mathbf{x}) - \gamma \quad (234)$$

The general form of the max entropy dist then given by

$$p(\mathbf{x}) = \frac{1}{Z} e^{\sum_i \gamma_i f_i(\mathbf{x})} \quad (235)$$

where $Z(\gamma_i) = \int d\mathbf{x} e^{\sum_i \gamma_i f_i(\mathbf{x})}$ = the partition function.

- The max ent dist = clearly just the usual Boltzmann dist, with energy $E(\mathbf{x}) = -\sum_i \gamma_i f_i(\mathbf{x})$.
- The values of the Lagrange multipliers are chosen to match the observed averages for the set of funcs $\{f_i(\mathbf{x})\}$ whose average value is being fixed:

$$\langle f_i \rangle_{model} = \int d\mathbf{x} p(\mathbf{x}) f_i(\mathbf{x}) = \frac{\partial \log Z}{\partial \lambda_i} = \langle f_i \rangle_{obs} \quad (236)$$

aka, the parameters of the dist can be chosen such that

$$\partial_{\lambda_i} \log Z = \langle f_i \rangle_{data} \quad (237)$$

- To gain more intuition for the MaxEnt dist, helpful to relate the Lagrange multipliers to the familiar thermodynamic quantities we use to describe physical systems.

- Our \mathbf{x} = the microscopic state of the system = the MaxEnt dist is a prob dist over microscopic states.
- But, in thermodynamics we only have access to average quantities. If we know only the average energy $\langle E(\mathbf{x}) \rangle_{obs}$, the MaxEnt procedure tells us to maximize the ent subject to the average energy constraint. This yields

$$p(\mathbf{x}) = \frac{1}{Z} e^{-\beta E(\mathbf{x})} \quad (238)$$

where we have identified the Lagrange multiplier conjugate to the energy $\lambda_1 = -\beta = 1/k_B T$ with the (negative) inverse temperature.

- Suppose we also constrain the particle number $\langle N(\mathbf{x}) \rangle_{obs}$. Then, an almost identical calc yields a MaxEnt dist of the functional form

$$p(\mathbf{x}) = \frac{1}{Z} e^{-\beta(E(\mathbf{x}) - \mu N(\mathbf{x}))} \quad (239)$$

where we have written our Lagrange multipliers in the familiar thermodynamic notation $\lambda_1 = -\beta$ and $\lambda_2 = \mu/\beta$.

- Since this is just the Boltzmann dist, we can also relate the partition func in our MaxEnt model to the thermodynamic free-energy via $F = -\beta^{-1}\log Z$.
- The choice of which quantities to constrain = equivalent to working in different thermodynamic ensembles.

15.2.2 From statistical mechanics to machine learning

- The MaxEnt idea also provides a general procedure for learning a generative model from *data*. Key difference between MaxEnt models in (theoretical) physics and ML:
 - In ML we have no direct access to observed average values $\langle f_i \rangle_{obs}$.
 - Instead, these averages must be directly estimated from data (samples).
 - To denote this difference, will call empirical averages calc from data as $\langle f_i \rangle_{data}$.

Can think of MaxEnt as a stat inference procedure simply by replacing $\langle f_i \rangle_{obs}$ by $\langle f_i \rangle_{data}$ above.

- This subtle change has important implications for training MaxEnt models.
 1. Since we don't know these averages exactly, but must estimate from the data, our training procedure must be careful not to overfit to the observations (our samples might not be reflective of the true values of these statistics).
 2. The averages of certain functions f_i are easier to estimate from limited data than others. Often an important consideration when formulating which MaxEnt model to fit to the data.
 3. Unlike in physics where conservation laws often guide the funcs f_i whose averages we hold fixed, ML offers no comparable guide for how to choose the f_i we care about.

For these reasons, choosing the $\{f_i\}$ is often far from straightforward.

- We have presented a physics based perspective of the MaxEnt procedure. The MaxEnt in ML is also closely related to Bayesian inference ideas and this latter point of view is more common in discussions in the stat and ML literature.

15.2.3 Generalized Ising Models from MaxEnt

- Form of a MaxEnt = completely specified once we choose the averages $\{f_i\}$ we wish to constrain.
- A common choice = constrain the first two moments of a dist.
- When our random variables \mathbf{x} =continuous, the corresponding MaxEnt dist = a multi-dim Gaussian.
- If \mathbf{x} =binary (discrete), corresponding MaxEnt= a generalized Ising (Potts) model with all-to-all couplings.
- To see this, consider
 - \mathbf{x} = a random variable with first and second moments $\langle x_i \rangle_{data}$ and $\langle x_i x_j \rangle_{data}$, respectively.
 - According to the Principle of Max Ent, we should choose to model this variable using a Boltzmann dist w/constraints on the first and second moments.
 - Let a_i be the Langrange multiplier associated with $\langle x_i \rangle_{data}$ and $J_{ij}/2$ be the L multipilier associated with $\langle x_i x_j \rangle_{data}$.
 - Using equation 236, it's easy to verify that the energy function

$$E(\mathbf{x}) = - \sum_i a_i x_i - \frac{1}{2} \sum_{ij} J_{ij} x_i x_j \quad (240)$$

satisfies the above constraints.

- Partition funcns for MaxEnt models=often intractable to compute. \rightarrow helpful to consider two special cases where \mathbf{x} has different support (different kinds of data).
- 1. Consider the case that the random variables $\mathbf{x} \in \mathbb{R}^n$ = real numbers. In this case, can compute the partition funcn directly:

$$Z = \int d\mathbf{x} e^{\mathbf{a}^T \mathbf{x} + \frac{1}{2} \mathbf{x}^T \mathbf{J} \mathbf{x}} = \sqrt{(2\pi)^n \det \mathbf{J}^{-1}} e^{-\frac{1}{2} \mathbf{a}^T \mathbf{J}^{-1} \mathbf{a}} \quad (241)$$

The resulting prob density funcn is

$$\begin{aligned} p(\mathbf{x}) &= Z^{-1} e^{-E(\mathbf{x})} \\ &= \frac{1}{\sqrt{(2\pi)^n \det \mathbf{J}^{-1}}} e^{-\frac{1}{2} \mathbf{a}^T \mathbf{J}^{-1} \mathbf{a} + \mathbf{a}^T \mathbf{x} + \frac{1}{2} \mathbf{x}^T \mathbf{J} \mathbf{x}} \\ &= \frac{1}{\sqrt{(2\pi)^n \det \Sigma}} e^{-\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu})} \end{aligned} \quad (242)$$

where $\boldsymbol{\mu} = -\mathbf{J}^{-1} \mathbf{a}$ and $\Sigma = -\mathbf{J}^{-1}$. This, ofc = the normalized, multi-dimensional Gaussian distribution.

2. Consider the case that the random variable \mathbf{x} = binary with $x_i \in \{-1, +1\}$. The energy func takes same form as the one above, but partition function can no longer be computed in a closed form. This model=the **Ising model** in physics and a **Markov Random Field** in ML. Since energy func intractable, the best we can do is estimate it using numerical techniques such as MCMC methods or approximate methods like variational MFT methods. Note: in ML common to use binary variables which take on values in $x_i \in \{0, 1\}$, rather than $x_i \in \{-1, +1\}$. Can sometimes be a source of confusion when translating between ML and physics literatures and can lead to confusion when using ML packages for physics problems.

15.3 Cost functions for training energy-based models

- MaxEnt procedures gives us a way of parameterizing an energy-based generative model.
- **For any energy-based generative model**, the energy function $E(\mathbf{x}, \{\theta_i\})$ depends on some parameters θ_i (=couplings in stat phys language) that must be inferred directly from the data.
 - Fex: For MaxEnt models the $\{\theta_i\}$ =just the Langrange multipliers $\{\lambda_i\}$ introduced in the last section.

The goal of the training procedure is to use the available data to fit these parameters.

- We fit the parameters by minimizing a cost func using stochastic gradient descent. The procedure naturally separates into two parts: choosing an appropriate cost func and calc the gradient of the cost func wrt the model parameters.
- Formulating a cost func for generative models = a little bit trickier than for supervised, discriminative models. Objective of discriminative models=straightforward: predict the label from the features. What we mean by a "good" generative model =much harder to define using a cost func.
 - Would like the model to generate examples similar to those found in the training data
 - But would also like it to be able to generalize - we don't want it to reproduce "spurious details" that are particular to the training data

- **Unlike for discriminative models, no straightforward idea like cross-validation on the data labels** that neatly addresses this issue.
- Calc the gradients of energy-based models also different from discriminative models such as deep NNs. Rather than rely on automatic differentiation techniques and backprop, calc the gradient requires drawing on intuitions from MCMC methods.
- Below provide an in-depth discussion of Boltzmann learning for energy-based generative models, focusing on MaxEnt models. Put emphasis on training procedures that generalize to more complicated generative models w/latent variables such as RBMs. Thus, largely ignore the incredibly rich physics-based literature on fitting Ising-like MaxEnt models.

15.3.1 Maximum likelihood

- By far most common approach for training a generative model = maximize the log-likelihood of the training data.
- Recall: the log-likelihood characterizes the log-prob of generating observed data using our generative model. By choosing the negative log-likelihood as the cost func, the learning procedure tries to find parameters that maximize the prob of the data.
- This cost func=intuitive, has therefore been the work-horse of most generative modeling. But, note the MLE has some important limitations we'll return to in chapter 12.
- We now employ a general notation applicable to all energy-based models, not just the MaxEnt. Reason being much of the discussion does not rely on the specific form of the energy func,, only the fact that our generative model takes a Boltzmann form.
 - $p_{\theta}(\mathbf{x})$ = the prob dist by which we denote the generative model. $\log Z(\{\theta_i\})$ = its corresponding partition func.
 - In MLE, the parameters are fit by maximizing the log-likelihood:

$$\mathcal{L}(\{\theta_i\}) := \langle \log(p_{\theta}(\mathbf{x})) \rangle_{data} \quad (243)$$

$$= - \langle E(\mathbf{x}; \{\theta_i\}) \rangle_{data} - \log Z(\{\theta_i\}) \quad (244)$$

where we set $\beta = 1$. Made use of the fact that our generative dist is of the Boltzmann form and that our partition func does not depend on the dat:

$$\langle \log Z(\{\theta_i\}) \rangle_{data} = \log Z(\{\theta_i\}) \quad (245)$$

15.3.2 Regularization

- Just as for discriminative models like linear and logistic regression, it's common to supplement the log-likelihood w/additional regularization terms that prevent overfitting.
- Instead of minimizing the log-likelihood, minimize a cost func of the form

$$-\mathcal{L}(\{\theta_i\}) + E_{reg}(\{\theta_i\}) \quad (246)$$

where $E_{reg}(\{\theta_i\})$ = an additional regularization term that prevents overfitting.

- **From Bayesian perspective**, this new term can be viewed as encoding a (negative) log-prior on model parameters and performing a maximum-a-posteriori (MAP) estimate instead of a MLE.
- As seen with regression, different forms of regularization give rise to different kinds of properties. A common choice for the regularization func are the sums of the L_1 or L_2 norms of the parameters

$$E_{reg}(\{\theta_i\}) = \Lambda \sum_i |\theta_i|^\alpha, \quad \alpha = 1, 2 \quad (247)$$

with Λ =parameters that control the regularization strenght. $\Lambda = 0$ =no regularization, simply performing MLE. Large Λ =will force many parameters to be close to or exactly zero.

- Just as in regression, an L_1 penalty enforces sparsity, w/many of the θ_i set to zero, and L_2 regularization shrinks the size of the parameters towards zero.
- One challenge of generative models: often hard to choose the reg strength Λ .
 - Recall: for linear and logistic regression Λ is chosen to maximize out-of-sample performance on a validation dataset (i.e. cross-validation).
 - But, for generative models data usually unlabeled. \rightarrow choosing Λ more subtle, no universal procedure.
 - One common strategy: divide the data into a training set and a validation set and monitor a summary statistic such as the
 - * log-likelihood
 - * energy distance (ref)
 - * variational free-energy of the generative model

on the training and validation sets. If the gap between them starts growing, one is probably overfitting the model even if the log-likelihood of the training dataset is still increasing.

- This also gives a procedure for "early stopping" - a regularization procedure we introduced in the context of discriminative models.
- In practice when using such regularizers: important to try many different values of Λ and then try to use a proxy statistic for overfitting to evaluate the optimal choice of Λ .

15.4 Computing gradients

- Now: procedure for minimizing the cost func. Powerful and common choice = SGD.
- Performing MLE using SGD requires calc the gradient of the log-likelihood wrt the parameters θ_i .
- To simplify notation and gain intuition, helpful to define "operators" $O_i(\mathbf{x})$, conjugate to the parameters θ_i

$$O_i(\mathbf{x}) := \frac{\partial E(\mathbf{x}; \theta_i)}{\partial \theta_i} \quad (248)$$

- Since the partition func just the cumulant generative func for the Boltzmann dist, we know the usual stat mech relationships between expectation values and derivatives of the log-partition func hold:

$$\langle O_i(\mathbf{x}) \rangle_{model} = \text{Tr}_{\mathbf{x}} p_{\theta}(\mathbf{x}) O_i(\mathbf{x}) = - \frac{\partial \log Z(\{\theta_i\})}{\partial \theta_i} \quad (249)$$

- In terms of the operators $\{O_i(\mathbf{x})\}$, the gradient of the log-likelihood takes the form

$$\begin{aligned} - \frac{\partial \mathcal{L}(\{\theta_i\})}{\partial \theta_i} &= \left\langle \frac{\partial E(\mathbf{x}; \theta_i)}{\partial \theta_i} \right\rangle_{data} + \frac{\partial \log Z(\{\theta_i\})}{\partial \theta_i} \\ &= \langle O_i(\mathbf{x}) \rangle_{data} - \langle O_i(\mathbf{x}) \rangle_{model} \end{aligned} \quad (250)$$

These eqs have a simple and beautiful interpretation:

- The grad of the log-likelihood wrt to a model parameter is a **difference of moments** - one calc directly from the data, one calc from our model using the current model parameters.
- *The positive phase* of the gradient = the data dependent term.
- *The negative phase* of the gradient = the model dependent term.

This derivation also gives an intuitive explanation for likelihood-based training procedures: The gradient acts on the model to lower the energy of configurations that are near observed data points while raising the energy of configs far from observed data points.

- Note: All info about the data only enters the training procedure through the expectations $\langle O_i(\mathbf{x}) \rangle_{data}$ and our generative model is blind to information beyond what is contained in these expectations.
- Now still need to calc the expectation values in the above eq. The positive gradient phase/ expec vals wrt the data is easily calc using samples from the training data. But, the negative phase/the expec vals wrt the model generally much harder to compute. In almost all cases will have to resort to numerical or approximate methods. Fundamental reason: impossible to calc the partition func exactly for most interesting problems.
- There are exceptional cases where we can calc expec vals analytically. The generative model then said to have a *Tractable Likelihood*. One example:
 - The Gaussian MaxEnt model for real valued data discussed before. The parameters/Lagrange multipliers for this model = the local fields \mathbf{a} and the pairwise coupling matrix J .
 - Here, the usual manipulations involving Gaussian integrals allow us to exactly find the parameters $\mu = -J^{-1}\mathbf{a}$ and $\Sigma = -J^{-1}$, yielding the familiar expressions $\mu = \langle \mathbf{x} \rangle_{data}$ and $\Sigma = \langle (\mathbf{x} - \langle \mathbf{x} \rangle_{data})(\mathbf{x} - \langle \mathbf{x} \rangle_{data})^T \rangle_{data}$. These are the standard estimates for the sample mean and covariance matrix.
 - Converting back to the Lagrange multiplier yields

$$J = -\langle (\mathbf{x} - \langle \mathbf{x} \rangle_{data})(\mathbf{x} - \langle \mathbf{x} \rangle_{data})^T \rangle_{data}^{-1} \quad (251)$$

SPM: men vil dette være overfittet...??? hvis det er perfekt til dataene? eller tar likelihood og skiller perfekt på hvor denne grensen går? hvordan?

- In the generic case with *intractable likelihoods*, we must estimate expec vals numerically. One way to do this:
 - Draw samples $\mathcal{S}_{model} = \{\mathbf{x}'_i\}$ from the model $p_\theta(\mathbf{x})$ and evaluate expec vals using these samples:

$$\langle h(\mathbf{x}) \rangle_{model} = \int d\mathbf{x} p_\theta(\mathbf{x}) h(\mathbf{x}) \approx \sum_{\mathbf{x}'_i \in \mathcal{S}_{model}} g(\mathbf{x}'_i) \quad (252)$$

The samples from the model $\mathbf{x}'_i \in \mathcal{S}_{model}$ are often referred to as *fantasy particles* in the ML literature and can be generated using simple MCMC algos such as Metropolis-Hastings.

- Once we have the fantasy particles (samples) from the model, we can also easily calc the grad of an arbitrary expec val $\langle g(\mathbf{x}) \rangle_{model}$ using what is commonly called the "log-derivative trick" in ML:

$$\frac{\partial}{\partial \theta_i} \langle g(\mathbf{x}) \rangle_{model} = \int d\mathbf{x} \frac{\partial p_\theta(\mathbf{x})}{\partial \theta_i} g(\mathbf{x}) \quad (253)$$

$$= \left\langle \frac{\partial \log p_\theta(\mathbf{x})}{\partial \theta_i} \right\rangle_{model} \quad (254)$$

$$= \langle O_i(\mathbf{x}) g(\mathbf{x}) \rangle_{model} \quad (255)$$

$$\approx \sum_{\mathbf{x}'_i \in \mathcal{S}_{model}} O_i(\mathbf{x}) g(\mathbf{x}'_i) \quad (256)$$

This expression allows us to take gradients of more complex cost funcs beyond the MLE procedure discussed here.

15.5 Summary of the training procedure

- Now summarize and present a general procedure for training an energy based model using SGD o the cost func.
- Goal: fit the parameters of a model $p_\lambda(\{\theta_i\}) = Z^{-1} e^{-E(\mathbf{x}, \{\theta_i\})}$
- Training the model involves the following steps
 1. Read a minibatch of data, $\{\mathbf{x}\}$
 2. Generate a random sample (fantasy particles) $\{\mathbf{x}'\} \sim p_\lambda$ using an MCMC algo (e.g. Metropolis-Hastings)
 3. Compute the grad of log-likelihood using these samples and the eq found earlier, where the averages are taken over the minibatch of data and the fantasy particles/samples from the model, respectively.
 4. Use the grad as input to one of the grad based optimizers discussed in chapter four.
- In practice: helpful to supplement this with some practical tricks that help training. As with discriminative NNs, important to initialize the parameters properly and print summary statistics during the training on the training and validation sets to prevent overfitting. **These and many other little practical tricks have been nicely summarized in a short note from the Hinton group REF**

- A major computational and practical limitation: often hard to draw samples from generative models. MCMC methods often have long mixing-times (=the time you have to run the Markov chain to get uncorrelated samples) → can result in biased sampling. Luckily, often don't need to know the gradients exactly for training ML models (recall: noisy gradient estimates often help the convergence of gradient descent algos) → can significantly reduce the computational expense by running MCMC for a reasonable time window. Will exploit this extensively in the next section discussing how to train more complex energy-based models with hidden variables.

16 Deep generative models: Latent variables and Restricted Boltzmann Machines (RBMs)

- Last section: core ideas behind energy-based generative models. Now: energy-based models that include latent/hidden variables.
- Including latent variables = greatly enhances expressive power, allows model to represent sophisticated correlations between visible features without sacrificing trainability. Multiple latent layers: can construct powerful deep generative models that possesses many of the same desirable properties as the deep, discriminative NNs.

16.1 Why hidden (latent) variables?

- Latent variables = powerful yet elegant way to encode sophisticated correlations between observable features. Underlying reason:
 - Marginalizing over a subset of variables (= "integrating out" degrees of freedom in physics language) induces complex interactions between remaining variables.
 - The idea that integrating out variables can lead to complex correlations = familiar component of many physical theories. Fex: when considering free electrons living on a lattice, integrating out phonons gives rise to higher-order electron-electron interactions (e.g. superconducting or magnetic correlations).
 - More generally, in the Wilsonian renormalization group paradigm, all effective field theories can be thought of as arising from integrating out high-energy degrees of freedom **REF**.
- Generative models w/latent variables run this logic in reverse:
 - Encode complex interactions between visible variables by introducing additional, hidden variables that interact w/visible degrees

of freedom in a simple manner, yet still reproduce the complex correlations between visible degrees in the data once marginalized over (integrated out).

- This trick is **also widely exploited in physics**, e.g. in the Hubbard-Stratonovich transformation or in the introduction of ghost fields in gauge theory.
- To make these ideas more concrete, let's revisit the pairwise Ising model introduced in discussion of MaxEnt models:

- The model is described by a Boltzmann distribution w/energy

$$E(\mathbf{v}) = - \sum_i a_i v_i - \frac{1}{2} \sum_{ij} v_i J_{ij} v_j \quad (257)$$

where J_{ij} = a symmetric coupling matrix that encodes the pairwise constraints and a_i enforce the single-variable constraint (single-variable constraint?).

- Goal: replace the complicated interactions between the visible variables v_i , encoded by J_{ij} , by interactions with a new set of latent variables h_μ .
- In order to do this, helpful to rewrite the coupling matrix in a slightly different form. Using SVD, we can always express the coupling matrix $J_{ij} = \sum_{\mu=1}^N W_{i\mu} W_{j\mu}$, where $\{W_{i\mu}\}_i$ = appropriately normalized singular vectors.
- In terms of $W_{i\mu}$, the energy takes the form

$$E_{Hop}(\mathbf{v}) = - \sum_i a_i v_i - \frac{1}{2} \sum_{ij\mu} v_i W_{i\mu} W_{j\mu} v_j \quad (258)$$

- Note: in the special case when both the $v_i \in \{-1, +1\}$ and $W_{i\mu} \in \{-1, +1\}$ are binary variables, a model with this form of the energy function is known as the **Hopfield model REF**. The Hopfield model has played an extremely important role in stat phys, computational neuroscience, and ML, see **REF** for a beautiful discussion combining all these properties.
- We therefore refer to all energy functions of the form above as **(generalized) Hopfield models**, even for the case when the $W_{i\mu}$ =continuous variables.
- Now "decouple" the v_i by introducing a set of normally, distributed continuous latent variables h_μ (**in condensed matter physics this called a Hubbard-Stratonovich transformation** as mentioned above).

- Using the usual identity for Gaussian integrals, we can rewrite the Boltzmann dist for the generalized Hopfield model as (**cool!! this could be used to show how to go from wave func sqaured to something involving hidden variables**)

$$p(\mathbf{v}) = \frac{e^{\sum_i a_i v_i + \frac{1}{2} \sum_{ij\mu} v_i W_{i\mu} W_{j\mu} v_j}}{Z} \quad (259)$$

$$= \frac{e^{\sum_i a_i v_i} \prod_{\mu} \int d h_{\mu} e^{-\frac{1}{2} \sum_{\mu} h_{\mu}^2 - \sum_i v_i W_{i\mu} h_{\mu}}}{Z} \quad (260)$$

$$= \frac{\int d\mathbf{h} e^{-E(\mathbf{v}, \mathbf{h})}}{Z} \quad (261)$$

where $E(\mathbf{v}, \mathbf{h})$ is a **joint energy functional** of both the latent and visible variables of the form

$$E(\mathbf{v}, \mathbf{h}) = - \sum_i a_i v_i + \frac{1}{2} \sum_{\mu} h_{\mu}^2 - \sum_{i\mu} v_i W_{i\mu} h_{\mu} \quad (262)$$

- We can also use the energy func $E(\mathbf{v}, \mathbf{h})$ to define a new energy-based model $p(\mathbf{v}, \mathbf{h})$ on both the latent and visible variables

$$p(\mathbf{v}, \mathbf{h}) = \frac{e^{-E(\mathbf{v}, \mathbf{h})}}{Z'} \quad (263)$$

- Marginalizing over latent variables of course gives us back the generalized Hopfield model (**REF**).

$$p(\mathbf{v}) = \int d\mathbf{h} p(\mathbf{v}, \mathbf{h}) = \frac{e^{-E_{Hop}(\mathbf{v})}}{Z} \quad (264)$$

- Notice: $E(\mathbf{v}, \mathbf{h})$ contains no direct interactions between visible degrees of freedom (or between hidden degrees of freedom). Instead: the complex correlations between the v_i are encoded in the interaction between the visible v_i and latent variables h_{μ} .
- It turns out the model presented here = a special case of a more general class of powerful energy-based models called Restrected Boltzmann Machines (RBMs).

16.2 Restricted Boltzmann machines (RBMs)

- RBM = an energy-based model w/both visible and hidden units where the visible and hidden units interact with each other but don't interact among themselves.

- Energy func of an RBM takes the general functional form

$$E(\mathbf{v}, \mathbf{h}) = - \sum_i a_i(v_i) - \sum_\mu b_\mu(h_\mu) - \sum_{i\mu} W_{i\mu} v_i h_\mu \quad (265)$$

where $a_i(\cdot)$ and $b_\mu(\cdot)$ are funcs that we are free to choose. The most common choice is:

$$a_i(v_i) := \begin{cases} a_i v_i, & \text{if } v_i \in \{0, 1\} \text{ is binary} \\ \frac{v_i^2}{2\sigma_i^2} & \text{if } v_i \in \mathbb{R} \text{ is continuous} \end{cases} \quad (266)$$

and

$$b_\mu(h_\mu) := \begin{cases} b_\mu h_\mu, & \text{if } h_\mu \in \{0, 1\} \text{ is binary} \\ \frac{h_\mu^2}{2\sigma_\mu^2} & \text{if } h_\mu \in \mathbb{R} \text{ is continuous} \end{cases} \quad (267)$$

For this choice of $a_i(\cdot)$ and $b_i(\cdot)$, layers consisting of discrete binary units often called Bernoulli layers and layers consisting of continuous variables often called Gaussian (WHAT happened to subtracting biases in the Gaussian version????).

- The **basic bipartite structure** of an RBM (=visible and hidden layer that interact with each other but not among themselves) often depicted using a graph as in fig 61.
- Different types:
 - Bernoulli-Bernoulli = most common choice
 - Continuous-continuous = **the RBM reduces to a multi-dimensional Gaussian with a very paritcular correlation structure.**
 - discrete (visible) - continuous (hidden) = the RBM is equivalent to a generalized Hopfield model.
 - continuous (visible) - discrete (hidden) = often called a Gaussian Bernoulli RBM (**REF, Hinton**).
 - Even possible to perform multi-modal learning w/mixture of cont and discrete variables.
 - Remember for all these: only interactions between hidden and visible, no intralayer interactions. **This is analogous to Quantum Electrodynamics**, where a free fermion and a free photon interact with one anotherbut not among themselves.
- Specifying a generative model w/this bipartite interaction structure has two major advantages:
 - Enables capturing both pairwise *and higher-order* correlations between visible units

- Makes it easier to sample from the model using an MCMC method known as block Gibbs sampling, in turn making the model easier to train.
- **It's worth better understanding the kind of correlations that can be captured using an RBM.**
 - To do so, can marginalize over the hidden units and ask about the resulting distribution over just the visible units

$$p(\mathbf{v}) = \int d\mathbf{h} p(\mathbf{v}, \mathbf{h}) = \int d\mathbf{h} \frac{e^{-E(\mathbf{v}, \mathbf{h})}}{Z} \quad (268)$$

(where the integral should be replaced by a trace in all expressions for discrete units).

- Can also define a marginal energy using the expression

$$p(\mathbf{v}) := \frac{e^{E(\mathbf{v})}}{Z} \quad (269)$$

- Combining these equations (how did you get rid of Z ?),

$$E(\mathbf{v}) = -\log \int d\mathbf{h} e^{-E(\mathbf{v}, \mathbf{h})} \quad (270)$$

$$= -\sum_i a_i(v_i) - \sum_\mu \log \int dh_\mu e^{b_\mu(h_\mu) + \sum_i v_i W_{i\mu} h_\mu} \quad (271)$$

- **To understand what correlations are captured by $p(\mathbf{v})$** it is helpful to introduce the distribution

$$q_\mu(h_\mu) = \frac{e^{b_\mu(h_\mu)}}{Z} \quad (272)$$

of hidden units h_μ , ignoring the interactions between \mathbf{v} and \mathbf{h} , and the **cumulant generating function** (???)

$$K_\mu(t) := \log \int dh_\mu q_\mu(h_\mu) e^{th_\mu} = \sum_n \kappa_\mu^{(n)} \frac{t^n}{n!} \quad (273)$$

$K_\mu(t)$ is defined such that the n^{th} cumulant is $\kappa_\mu^{(n)} = \partial_t^n K_\mu|_{t=0}$.

- The cumulant generating func appears in the **marginal free-energy of the visible units**, which can be rewritten (up to a

constant term) as:

$$E(\mathbf{v}) = - \sum_i a_i(v_i) - \sum_{\mu} K_{\mu} \left(\sum_i W_{i\mu} v_i \right) \quad (274)$$

$$= - \sum_i a_i(v_i) - \sum_{\mu} \sum_n \kappa_{\mu}^{(n)} \frac{(\sum_i W_{i\mu} v_i)^n}{n!} \quad (275)$$

$$= - \sum_i a_i(v_i) - \sum_i \left(\sum_{\mu} \kappa_{\mu}^{(1)} W_{i\mu} \right) v_i \quad (276)$$

$$- \frac{1}{2} \sum_{ij} \left(\sum_{\mu} \kappa_{\mu}^{(2)} W_{i\mu} W_{j\mu} \right) v_i v_j + \dots \quad (277)$$

We see that

- * The marginal energy includes all orders of interactions between the visible units, with the n -th order cumulants of $q_{\mu}(h_{\mu})$ weighting the n -th order interactions between the visible units. (!!!!!)
- * In the case of the Hopfield model we discussed previously, $q_{\mu}(h_{\mu})$ = a standard Gaussian distribution where the mean is $\kappa_{\mu}^{(1)} = 0$, the variance is $\kappa_{\mu}^{(2)} = 1$, and all higher order cumulants are zero. Plugging these cumulants (weren't these called moments earlier???? same thing or different?) into the above equation recovers the expression we first got for the energy $E(\mathbf{v}, \mathbf{h})$ after having introduced the hidden variables before (one section back).
- These calc make clear the underlying reason for the incredible representational power of RBMs with a Bernoulli hidden layer (did we assume the type of the hidden layer above?).
 - **Each hidden unit can encode very complex interactions at all orders.**
 - We can learn which order of correlations/interactions are important directly from the data instead of having to specify them ahead of time as we did in the MaxEnt models.
 - This highlights the power of generative models with even the simplest interactions between visible and latent variables to encode, learn, and represent complex correlations present in the data.

16.3 Training RBMs

- RBMs—a special case of energy-based generative models, which can be trained using the MLE procedure, as described in prev chapter. Brief recap:

- First, choose a cost func: For MLE, this just the negative log-likelihood with or without an additional regularization term to prevent overfitting.
- Then, minimize the cost func using one of the SGD methods. The gradient itself can be calc using the equation from prev chapter. Fex: for the Bernoulli-Bernoulli RBM we have

$$\frac{\partial \mathcal{L}(\{W_{i\mu}, a_i, b_\mu\})}{\partial W_{i\mu}} = \langle v_i h_\mu \rangle_{data} - \langle v_i h_\mu \rangle_{model} \quad (278)$$

$$\frac{\partial \mathcal{L}(\{W_{i\mu}, a_i, b_\mu\})}{\partial a_i} = \langle v_i \rangle_{data} - \langle v_i \rangle_{model} \quad (279)$$

$$\frac{\partial \mathcal{L}(\{W_{i\mu}, a_i, b_\mu\})}{\partial b_\mu} = \langle h_\mu \rangle_{data} - \langle h_\mu \rangle_{model} \quad (280)$$

$$(281)$$

where

- * The postive expec wrt data = sampling from the model while clamping the visible units to their observed values in the data.
- * As before, calc the negative phase of the gradient (the expec val wrt the model) requires we draw samples from the model. Luckily, the bipartite form of the interactions in RBMs were specifically chosen with this in mind.

16.3.1 Gibbs sampling and contrastive divergence (CD)

- The **bipartite** interaction structure of the RBM makes it possible to calc expec vals using a MCMC method called Gibbs sampling.
- Key reason for this: since no interactions of visible units with themselves or hidden with themselves; the visible and hidden units are conditionally independent:

$$p(\mathbf{v}|\mathbf{h}) = \prod_i p(v_i|\mathbf{h}) \quad (282)$$

$$p(\mathbf{h}|\mathbf{v}) = \prod_\mu p(h_\mu|\mathbf{v}) \quad (283)$$

with

$$p(v_i = 1|\mathbf{h}) = \sigma(a_i + \sum_\mu W_{i\mu} h_\mu) \quad (284)$$

$$p(h_\mu = 1|\mathbf{v}) = \sigma(b_\mu + \sum_i W_{i\mu} v_i) \quad (285)$$

and where $\sigma(x) = 1/(1 + e^{-x})$ = the sigmoid function.

- Using these expressions, easy to compute expect vals wrt the data.
 - Input to grad descent = a minibatch of observed data.
 - For each sample in the minibatch, we simply clamp the visible units to the observed vals and apply the above equation using the probability of the hidden variables.
 - Then average over all samples in the minibatch to calc expect vals wrt the data.
- To calc expect vals wrt the model, we use (block) Gibbs sampling. Idea:
 - Iteratively sample from the conditional distributions $\mathbf{h}_{t+1} \sim p(\mathbf{h}|\mathbf{v}_t)$ and $\mathbf{v}_{t+1} \sim p(\mathbf{v}|\mathbf{h}_{t+1})$.
 - Since the units conditionally independent, each step of this iteration can be performed by simply drawing random numbers.
 - The samples are guaranteed to converge to the equilibrium distribution of the model in the limit that $t \rightarrow \infty$.
 - At the end of the Gibbs sampling, one ends up with a minibatch of samples (fantasy particles).
- One drawback of Gibbs sampling: may take many back and forth iterations to draw an independent samples. \rightarrow the Hinton group introduced an approximate Gibbs sampling technique called Contrastive Divergence (CD).
 - In CD- n , we just perform n iterations of (block) Gibbs sampling, with n often taken to be as small as 1.
 - Price for this truncation: we're not drawing samples from the true model distribution.
 - But for our purpose - using the expects to estimate the gradient for SGD - CD- n has been proven to work reasonably well. As long as the approximate gradients are reasonably correlated with the true gradient, SGD will move in a reasonable direction.
 - CD- n of course does come at a price. Truncating the Gibbs sampler prevents sampling far away from the starting point, which for CD- n are the data points in the minibatch. \rightarrow our generative model will be much more accurate around regions of feature space close to our training data. \rightarrow as is often the case in ML, CD- n sacrifices ability to generalize to some extent in order to make the model easier to train.
 - Some of these undesirable features can be tempered by using a slightly different variant of CD called Persistent Contrastive Divergence (PCD).

- * Rather than restarting the Gibbs sampler from the data at each gradient descent step, we start the Gibbs sampling at the fantasy particles (samples from the model) in the last grad descent step.
- * Since parameters change slowly compared to the Gibbs sampling, samples that are high prob at one step of the SGD are also likely to be high prob at the next step → ensures PCD does not introduce large errors in the estimation of the gradients.
- * Advantage of using fantasy particles to initialize the Gibbs sampler = to allow PCD to explore parts of the feature space much further from the training data than one could reach with ordinary CD.

16.3.2 Practical considerations

"Tricks-of-the-trade" training of RBM summary published in note by Hinton **REF**. Some of the important points:

- **Initialization** The model must be initialized.
 - Hinton suggests taking the weights $W_{i\mu}$ from a Gaussian with a mean zero and std $\sigma = 0.01$.
 - Alternative, proposed by Glorot and Bengio: choose std to scale with the layer sizes: $\sigma = 2/\sqrt{N_v + N_h}$ where N_v and N_h =number of visible and hidden units respectively.
 - Bias of the hidden units initialized to zero
 - Bias of the visible units typically taken to be inversely proportional to the mean activation, $a_i = \langle v_i \rangle_{data}^{-1}$
- **Regularization**
 - Can use an L_1 or L_2 penalty, typically only on the weight parameters, not the biases.
 - Alternatively, Dropout has been shown to decrease overfitting when training with CD and PCD and to result in more interpretable learned features.
- **Learning Rates** Typically, helpful to reduce the learning rate in later stages of training.
- **Updates for CD and PCD** Several computational tricks can be used for speeding up the alternating updates in CD and PCD.

16.4 Deep Boltzmann Machine

- Possess multiple hidden layers and were the first models rebranded as "deep learning" (technically, these were Deep Belief Networks where only the top layer was undirected).
- Motivation:
 - An RBM is composed of two layers of neurons, connected via an **undirected graph**. As a result, possible to perform sampling $\mathbf{v} \sim p(\mathbf{v}|\mathbf{h})$ and inference $\mathbf{h} \sim p(\mathbf{h}|\mathbf{v})$ w/the same model.
 - As w/the Hopfield model, can view **each hidden unit as representative of a pattern, or feature, that could be present in the data (in general, should think of activity patterns of hidden units representing features in the data)**. The inference step involves assigning a prob to each of these features that expresses **the degree to which each feature is present in the data sample**.
 - In an RBM, hidden units don't influence each other during the inference step, i.e. hidden units are conditionally independent given the visible units. There are a number of reasons why this is unsatisfactory. One is:
 - * The desire for sparse, distributed representations, where each observed visible vector will strongly activate a few (i.e. more than one but only a very small fraction) of the hidden units.

In the brain, this is thought to be achieved by inhibitory lateral connections between neurons. But, adding lateral intra-layer connections between the hidden units makes the dist difficult to sample from so we need to come up with another way of creating connections between the hidden units.

 - With the Hopfield model, we saw that pairwise linear connections between neurons can be mediated through another layer. \rightarrow a simple way to allow for effective connections between the hidden units is to add another layer of hidden units.
 - Rather than just having two layers, a visible and hidden, we can add additional layers of latent variables to account for the correlations between hidden units.
 - Ideally, as one adds more and more layers, one might hope that the correlations between hidden variables become smaller and smaller deeper into the network. This basic logic is **reminiscent of renormalization procedures that seek to decorrelate layers at each step REF**.
 - Price of adding additional layers = models become harder to train.

- Training DBMs more subtle than RBMs due to the difficulty of propagating information from visible to hidden units. But, Hinton and co realized some of these problems could be alleviated via a layerwise procedure.
 - Rather than attempting to train the whole DBM at once, can think of the DBM as a stack of RBMs (see fig 63)
 - One first trains the bottom two layers of the DBM - treating it as if it is a standalone RBM.
 - Once bottom trained, can generate "samples" from the hidden layer and use these samples as an input to the next RBM (consisting of the first and second hidden layer - purple hexagons and green squares in fig 63).
 - This procedure can be repeated to pretrain all layers of the DBM.
 - This pretraining initializes the weights so that SGD can be used effectively when the network is trained in a supervised (???) fashion.
 - In particular, the pretraining helps the grads to stay well behaved rather than vanish or blow up - a problem we discussed extensively in the section on DNNs.
 - Worth noting: once pretrained we can use the usual Boltzmann learning rules (presented in the prev chapter) to fine-tune the weights and improve the performance of the DBM **REF**.
 - As shown in next section, Paysage package presented here can be used to both construct and train DBNs using such a pretraining procedure.

16.5 Example: Using Paysage for MNIST

- Here demonstrate how to use the new open source package Paysage (French for landscape) for training unsupervised energy-based models on the MNIST dataset. (Paysage documentation link to github).
- We'll show how to build and train four different models:
 - A "Hopfield" type RBM with Gaussian hidden and Bernoulli (binary) visible units.
 - A conventional Bernoulli-Bernoulli RBM
 - A conventional RBM with an additional L_1 -penalty that enforces sparsity
 - A DBM with three Bernoulli layers with L_1 penalty each.

- Note Paysage requires Python 3.6 or higher. We fix the seed of the rand num generator to ensure reproducibility.
- Procedure:
 - Download a preprocessed version of the MNIST data built into Paysage
 - If the first time using it, need to shuffle it. Necessary since we shall employ SGD-based algos which requires using small minibatches of data to compute the grad at each step. If the data ordered, then the estimates for the grads computed from the minibatches will be biased. Shuffling ensures the grad estimates unbiased (though still noisy).
 - Create a python generator, which splits the `data` into a training and validation sets, and separates them into minibatches of size `batch_size`. Before we begin training, we set `data` to training mode.
 - To monitor progress of performance metrics during training, define the variable `performance` which tells Paysage to measure the reconstruction error from the validation set. **Possible metrics include**
 - * The reconstruction error (used in this example)
 - * Metrics related to difference in energy of random samples from the model
 - * See `metrics.md` in Paysage documentation for a complete list.
 - Now move on to construct a `hopfield` model.
 - * Use the `Model` class and with a visible `BernoulliLayer` and a hidden `GaussianLayer`.
 - * Also standardize the mean and variance of the Gaussian layer setting them to zero and unity, respectively (the nomenclature of Paysage here is inspired by the terminology in Variational Autoencoders).
 - * We chose to train the model with the `Adam` optimizer. To ensure convergence, attenuate the `learning_rate` hyperparameter according to a `PowerLawDecay` schedule: `learning_rate(t) = initial/(1 + coefficient×t)`. It will prove convenient to define the function `Adam_optimizer` for this purpose.
 - * Next, have to create the model. First, initialize the `model` using the `initialize` function attribute which accepts the `data` as a required argument. We choose the initialization routine `glorot` (as mentioned in the practical tricks section in this chapter).

- * Second, define an optimizer calling the func `Adam_optimizer`, and store the object under the name `opt`.
 - * To create an MCMC `sampler`, we use the method `from_batch` of the `SequentialMC` class, passing the `model` and the `data`.
 - * Last, create an `SGD` object called `trainer` to train the model using Persistent Contrastive Divergence (`pcd`) with a fixed number of `monte_carlo_steps`.
 - * We can also `monitor` the reconstruction error during training.
 - * Last, we train the model in epochs (cf. variable `num_epochs`), calling the `train` method of `trainer`.
 - * These steps=universal for shallow generative `models`, and it's convenient to combine them in the func `train_model`, which we shall use repeatedly.
- Can easily build a Bernoulli RBM and train it using the funcs defined above as follows. (shows code.)
 - Constructing a Bernoulli RBM with L_1 regularization also straightforward, using the `add_penalty` method which accepts a dictionary as input. Some layers may have multiple properties (such as the location and scale parameters of a Gaussian layer) so the dictionary key specifies to which property the penalty should be applied.
 - To define a DBM,
 - * Just add more layers, and an L_1 penalty for every layer.
 - * Recalling the essential trick with layer-wise pre-training to prepare the weights of the DBM, we define a `pretrainer` as an object of the `LayerwisePretrain` class (see code snippet below). This results in a slight modification of the function `train_model`.
- Having trained our models, let's see how they perform by computing some reconstructions and fantasy particles from the validation data.
 - Recall: a reconstruction \mathbf{v}' of a given data point \mathbf{x} is computed in two steps:
 - * Fix the visible layer $\mathbf{v} = \mathbf{x}$ to be the data, use MCMC sampling to find the state of the hidden layer \mathbf{h} which maximizes the prob dist $p(\mathbf{h}|\mathbf{v})$.
 - * Fixing the same obtained state \mathbf{h} , we find the reconstruction \mathbf{v}' of the original data point which maximizes the prob $p(\mathbf{v}'|\mathbf{h})$.
- in the case of a DBM, the forward pass continues until we reach at the last of the hidden layers, and the backward pass goes in reverse.

- A config sampled from an RBM needs to specify the values of both the visible and hidden units. Since the data only specify the visible units, we need to initialize some hidden values. The visible and hidden units stored in a `State` object.
 - To compute reconstructions, we define an MCMC `sampler` based on the trained `model`. The starting point for the MCMC sampler is set using the `set_state()` method.
 - To compute reconstructions, need to keep the prob dist learned by the generative `model` fixed which is done by the help of the `deterministic_iteration` function method, that takes in its first argument the number of passes (1 for a single $\mathbf{v} \rightarrow \mathbf{h} \rightarrow \mathbf{v}'$ pass), and the state of the sampler `sampler.state` as required arguments.
 - Can combine these steps in the func `compute_reconstructions`.
 - See fig 60 for results.
- Once we have the trained models ready, can use MCMC to draw samples from the corresponding prob dist, called "fantasy particles". To this end, let's
 - Draw a `random_sample` from the validation data and compute the `model_state`.
 - Next, we define an MCMC `sampler` based on the `model`, and set its state to `model_state`.
 - To compute fantasy particles, we do layer-wise Gibbs sampling for a total of `n_steps` equilibration steps.
 - The last step (controlled by the boolean `mean_field`) is a final mean-field iteration (see tricks discussed in **Hinton REF**).
 - Result in fig 64.
 - One can use generative models to reduce noise in images (**de-noising**). Let's randomly flip a fraction, `fraction_to_flip`, of the black&white bits in the validation data, and use the models defined above to reconstruct (de-noise) the digit images. Result in fig 65.

16.6 Example: Using Paysage for the Ising Model

- Can also use Paysage to analyze the 2D Ising data.
 - In prev sections, we used our knowledge of the critical point at $T_c/J \approx 2.26$ to label the spin configs and study the problem of classifying the states according to their phase of matter.

- But, in more complicated models, where the precise position of T_c is not known, one cannot label the states with such an accuracy, if at all.
 - As we explained, **generative models can be used to learn a variational approximation for the probability distribution that generated the data points.**
 - By using only the 2D spin configs, we now want to train a Bernoulli RBM, the fantasy particles of which are thermal Ising configs.
 - Unlike previous studies of the Ising dataset, here we perform the analysis at a fixed temperature T . We can then apply our model at three different values $T = 1.75, 2.25, 2.75$ in the ordered, critical and disordered regions, respectively.
 - Define a DBM with two hidden layers of N_{hidden} and $N_{hidden}/10$ units, respectively, and apply L_1 regularization to all weights.
 - As in the MNIST problem above, we apply layer-wise pre-training, and deploy Persistent Contrastive Divergence to train the DBM using ADAM.
 - One lesson from this problem—similar to real-life problems, this task is computationally intensive. The training time on present-day laptops easily exceeds that of previous studies from this review. → encourage the reader to try GPU-based training and study the resulting speed-up.
- See figs 66, 67, 68 for the results of the numerical experiment at $T/J = 1.75, 2.25, 2.75$ respectively, for a DBM with $N_{hidden} = 800$. Looking at the reconstructions and the fantasy particles, we see that our DBM works well in the disordered and critical regions. But, the chosen layer architecture not optimal for the ordered phase.

16.7 Generative models in physics

- Examples of generative models studied in physics:
 - In biophysics, dynamic Boltzmann dists have been used as effective models in chemical kinetics.
 - In stat phys, they were used to identify the criticality in the Ising model.
 - In parallel, tools from stat phys have been applied to analyze the learning ability of RBMs **REF**, characterizing the sparsity of the weights, the effective temperature, the nonlinearities in the activation functions of hidden units, and the adaptation of fields maintaining the activity in the visible layer (?? what does this last thing mean?).

- Spin glass theory motivated a deterministic framework for the training, evaluation, and use of RBMs;
- It was demonstrated that the training process in RBMs itself exhibits phase transitions
- Learning in RBMs was studied in the context of nonequilibrium thermodynamics and spectral dynamics
- Mean-field theory found application in analyzing DBMs **REF**
- Using generative models to improve Monte Carlo algos
- Ideas from quantum mechanics have been put forward to introduce improved speed-up in certain parts of the learning algos for Helmholtz machines
- **WATCH OUT VILDE:**
- Generative models have applications in the study of quantum systems too. Most notably, RBM-inspired variational asatzes were used to learn the prob dist associated with the absolute square of a quantum state **REF**, and
- In this context, RBMs are sometimes called Born machines **REF**.
- Further applications include the detection of order in low-energy product states, and
- Learning Einstein-Podolsky-Rosen correlations on an RBM.
- Inspired by the success of tensor networks in physics, the latter have been used as a basis for RBMs to extract the spatial geometry from entangelement **REF**, and
- Generative models based on matrix product states have been developed.
- Quantum entanglement was studied using RBM-encoded states **REF** and
- Tensor product based generative models have been used to understand MNIST and other ML datasets.

17 Variational AutoEncoders (VAEs) and Generative Adversarial Networks (GANs)

- We previously considered energy-based generative models. Here, we extend to two new generative model frameworks that have gained wide appeal:
 - Generative adversarial networks (GANs)
 - Variational autoencoders (VAEs)

- Unlike energy-based models, both these frameworks are based on differentiable neural networks and consequently can be trained using **back-prop** methods. VAEs in particular can be easily implemented and trained using high-level packages such as Keras making them an easy-to-deploy generative framework.
- They also differ from energy-based models in that they **don't directly seek to maximize likelihood**. Fex: GANs employ a novel cost func based on adversarial learning.
- VAEs and GANs are starting to make their way into physics (refs). More generally they've found important applications in many artistic and image manipulation tasks.

17.1 The limitations of maximizing Likelihood

- The Kullback-Leibler (KL)-divergence plays a central role in many generative models. Developing an intuition about KL-divergences = one of the keys to understanding why adversarial learning has proved such a powerful method for generative modeling.
- The KL divergence measures the similarity between two prob dists $p(\mathbf{x})$ and $q(\mathbf{x})$. Strictly speaking it's **not a metric because it is not symmetric and does not satisfy the triangle inequality**.
- Given two distributions, there are two distinct KL-divergences we can construct:

$$D_{KL}(p||q) = \int d\mathbf{x} p(\mathbf{x}) \log \frac{p(\mathbf{x})}{q(\mathbf{x})} \quad (286)$$

$$D_{KL}(q||p) = \int d\mathbf{x} q(\mathbf{x}) \log \frac{q(\mathbf{x})}{p(\mathbf{x})} \quad (287)$$

A related quantity called the Jensen-Shannon divergence,

$$JS(p, q) = \frac{1}{2} [D_{KL}(p||\frac{p+q}{2}) + D_{KL}(q||\frac{p+q}{2})] \quad (288)$$

does satisfy all the properties of a squared metric (i.e. the square root of the JS divergence is a metric).

- An important property of the KL-divergence we'll make use of repeatedly is its positivity: $D_{KL}(p||q) \geq 0$ with equality iff $p(\mathbf{x}) = q(\mathbf{x})$ almost everywhere.
- In generative models in ML, the two dists we're usually concerned with = the model dist $p_{\theta}(\mathbf{x})$ and the data dist $p_{data}(\mathbf{x})$.

- We of course want these models to be as similar as possible. But, **there are many subtelties about how we measure similarities that can have large consequences for the behaviour of training procedures.**
- Maximizing the log-likelihood of the data under the model = minimizing the KL divergence between the data dist and the model dist $D_{KL}(p_{data}||p_{\theta})$. To see this, we can rewrite the KL divergence as:

$$D_{KL}(p_{data}||p_{\theta}) = \int d\mathbf{x} p_{data}(\mathbf{x}) \log p_{data}(\mathbf{x}) \quad (289)$$

$$- \int d\mathbf{x} p_{data}(\mathbf{x}) \log p_{\theta}(\mathbf{x}) \quad (290)$$

$$= -S[p_{data}] - \langle \log p_{\theta}(\mathbf{x}) \rangle_{data} \quad (291)$$

Rearranging this eq, we have

$$\langle \log p_{\theta}(\mathbf{x}) \rangle_{data} = -S[p_{data}] - D_{KL}(p_{data}||p_{\theta}) \quad (292)$$

(hva er S her?? JS divergence?? No dont think so, thats what GANs use) The equivalence follows from the positivity of KL-divergence and the fact that the entropy of the data dist is constant.

- In contrast, the original formulation of GANs minimizes **an upper bound on the Jensen-Shannon divergence** between the model dist $p_{\theta}(\mathbf{x})$ and the data dist $p_{data}(\mathbf{x})$.
- This difference in objectives underlies the difference in behavior between GANs and likelihood based generative models. To see this, we can compare the behavior of the two KL-divergences $D_{KL}(p_{data}||p_{\theta})$ and $D_{KL}(p_{\theta}||p_{data})$. As illustrated in fig 69 and 70 (**useful**), though both of these KL-divergences measure similarities between the two dists, they are sensitive to very different things.
 - $D_{KL}(p_{\theta}||p_{data})$ is insensitive to setting $p_{\theta} \approx 0$ even when $p_{data} \neq 0$,
 - Whereas $D_{KL}(p_{data}||p_{\theta})$ punishes this harshly.
 - In contrast, $D_{KL}(p_{data}||p_{\theta})$ prefers models that have a high prob in regions with lots of training data points
 - Whereas $D_{KL}(p_{\theta}||p_{data})$ punishes models for putting high prob where there is no data.
- In the context of the above discussion, this suggests that the way likelihood-based methods are most likely to fail, is by improperly "filling in" any low-prob density regions between peaks in the data dist.

- In contrast, at least in principle, the Jensen-Shannon dist which underlies GANs is sensitive to both
 - Placing weight where there is data since it has information about $D_{KL}(p_{data}||p_{\theta})$ and
 - Not placing weight where no data has been observed (i.e. low-prob density regions) since it has info about $D_{KL}(p_{\theta}||p_{data})$.
- In practice:
 - $D_{KL}(p_{data}||p_{\theta})$ can be calc easily directly from the data using sampling
 - $D_{KL}(p_{\theta}||p_{data})$ is impossible to compute since we don't know $p_{data}(\mathbf{x})$. In particular, this integral cannot be calc using sampling since we cannot evaluate $p_{data}(\mathbf{x})$ at the locations of the fantasy particles.
- The idea of adversarial learning - circumnavigate this difficulty by using an adversarial learning procedure.
- Recall, $D_{KL}(p_{\theta}||p_{data})$ large when the model artificially over-weighs low-density regions near real peaks. Adversarial learning accomplishes this same task by teaching a discriminator network to distinguish between real data points and samples generated from the model.
- By punishing the model for generating points that can be easily discriminated from the data, adversarial learning decreases the weight of regions in the model space that are far away data points - regions that inevitably arise when maximizing likelihood.
- This core intuition implicitly underlies many adversarial training algos (though it has been recently suggested this may not be the entire story.)

- 17.2 Generative models and adversarial learning
- 17.3 Variational Autoencoders (VAEs)
 - 17.3.1 VAEs as variational models
 - 17.3.2 Training via the reparametrization trick
 - 17.3.3 Connection to the information bottleneck
- 17.4 VAE with Gaussian latent variables and Gaussian encoder
 - 17.4.1 Implementing the Gaussian VAE
 - 17.4.2 VAEs for the MNIST dataset
 - 17.4.3 VAEs for the 2D Ising model
- 18 Outlook
 - 18.1 Research at the intersection of physics and ML
 - 18.2 Topics not covered in review
 - 18.3 Rebranding Machine Learning as "Artificial Intelligence"
 - 18.4 Social Implications of Machine Learning