

A Kernel Theory of Modern Data Augmentation

Tri Dao¹, Albert Gu¹, Alexander J. Ratner¹, Virginia Smith², Christopher De Sa³, and
Christopher Ré¹

¹Department of Computer Science, Stanford University

²Department of Electrical and Computer Engineering, Carnegie Mellon University

³Department of Computer Science, Cornell University

{trid,albertgu,ajratner}@stanford.edu, smithv@cmu.edu, cdesa@cs.cornell.edu,
chrismre@cs.stanford.edu

March 21, 2019

Abstract

Data augmentation, a technique in which a training set is expanded with class-preserving transformations, is ubiquitous in modern machine learning pipelines. In this paper, we seek to establish a theoretical framework for understanding data augmentation. We approach this from two directions: First, we provide a general model of augmentation as a Markov process, and show that kernels appear naturally with respect to this model, even when we do not employ kernel classification. Next, we analyze more directly the effect of augmentation on kernel classifiers, showing that data augmentation can be approximated by first-order feature averaging and second-order variance regularization components. These frameworks both serve to illustrate the ways in which data augmentation affects the downstream learning model, and the resulting analyses provide novel connections between prior work in invariant kernels, tangent propagation, and robust optimization. Finally, we provide several proof-of-concept applications showing that our theory can be useful for accelerating machine learning workflows, such as reducing the amount of computation needed to train using augmented data, and predicting the utility of a transformation prior to training.

1 Introduction

The process of augmenting a training dataset with synthetic examples has become a critical step in modern machine learning pipelines. The aim of data augmentation is to artificially create new training data by applying transformations, such as rotations or crops for images, to input data while preserving the class labels. This practice has many potential benefits: Data augmentation can encode prior knowledge about data or task-specific invariances, act as regularizer to make the resulting model more robust, and provide resources to data-hungry deep learning models. As a testament to its growing importance, the technique has been used to achieve nearly all state-of-the-art results in image recognition [4, 11, 13, 33], and is becoming a staple in many other areas as well [40, 22]. Learning augmentation policies alone can also boost the state-of-the-art performance in image classification tasks [31, 7].

Despite its ubiquity and importance to the learning process, data augmentation is typically performed in an ad-hoc manner with little understanding of the underlying theoretical principles. In the field of deep learning, for example, data augmentation is commonly understood to act as a regularizer by increasing the number of data points and constraining the model [12, 42]. However, even for simpler models, it is not well-understood how training on augmented data affects the learning process, the parameters, and the decision surface of the resulting model. This is exacerbated by the fact that data augmentation is performed in diverse ways in modern machine learning pipelines, for different tasks and domains, thus precluding a general model of transformation. Our results show that regularization is only part of the story.

In this paper, we aim to develop a theoretical understanding of data augmentation. First, in Section 3, we analyze data augmentation as a Markov process, in which augmentation is performed via a random sequence of transformations. This formulation closely matches how augmentation is often applied in practice. Surprisingly, we show that performing k -nearest neighbors with this model asymptotically results in a kernel classifier, where the kernel is a function of the base augmentations. These results demonstrate that kernels appear naturally with respect to data augmentation, regardless of the base model, and illustrate the effect of augmentation on the learned representation of the original data.

Motivated by the connection between data augmentation and kernels, in Section 4 we show that a kernel classifier on augmented data approximately decomposes into two components: (i) an averaged version of the transformed features, and (ii) a data-dependent variance regularization term. This suggests a more nuanced explanation of data augmentation—namely, that it improves generalization *both by inducing invariance and by reducing model complexity*. We validate the quality of our approximation empirically, and draw connections to other generalization-improving techniques, including recent work in invariant learning [43, 24, 30] and robust optimization [28].

Finally, in Section 5, to illustrate the utility of our theoretical understanding of augmentation, we explore promising practical applications, including: (i) developing a diagnostic to determine, prior to training, the importance of an augmentation; (ii) reducing training costs for kernel methods by allowing for augmentations to be applied directly to features—rather than the raw data—via a random Fourier features approach; and (iii) suggesting a heuristic for training neural networks to reduce computation while realizing most of the accuracy gain from augmentation.

2 Related Work

Data augmentation has long played an important role in machine learning. For many years it has been used, for example, in the form of *jittering* and *virtual examples* in the neural network and kernel methods literatures [35, 34, 8]. These methods aim to augment or modify the raw training data so that the learned model will be invariant to known transformations or perturbations. There has also been significant work in incorporating invariance directly into the model or training procedure, rather than by expanding the training set. One illustrative example is that of tangent propagation for neural networks [36, 37], which proposes a regularization penalty to enforce local invariance, and has been extended in several recent works [32, 9, 43]. However, while efforts have been made that loosely connect traditional data augmentation with these methods [21, 43], there has not been a rigorous study on how these sets of procedures relate in the context of modern models and transformations.

In this work, we make explicit the connection between augmentation and modifications to the model, and show that prior work on tangent propagation can be derived as a special case of our more general theoretical framework (Section 5). Moreover, we draw connections to recent work on invariant learning [24, 30] and robust optimization [28], illustrating that data augmentation not only affects the model by increasing invariance to specific transformations, but also by reducing the variance of the estimator. These analyses lead to an important insight into how invariance can be most effectively applied for kernel methods and deep learning architectures (Section 5), which we show can be used to reduce training computation and diagnose the effectiveness of various transformations.

Prior theory also does not capture the complex process by which data augmentation is often applied. For example, previous work [1, 3] shows that adding noise to input data has the effect of regularizing the model, but these effects have yet to be explored for more commonly applied complex transformations, and it is not well-understood how the inductive bias embedded in complex transformations manifest themselves in the invariance of the model (addressed here in Section 4). A common recipe in achieving state-of-the-art accuracy in image classification is to apply a sequence of more complex transformations such as crops, flips, or local affine transformations to the training data, with parameters drawn randomly from hand-tuned ranges [4, 10]. Similar strategies have also been employed in applications of classification for audio [40] and text [22]. In Section 3, we analyze a motivating model reaffirming the connection between augmentation and

kernel methods, even in the setting of complex and composed transformations.

Finally, while data augmentation has been well-studied in the kernels literature [2, 34, 25], it is typically explored in the context of simple geometrical invariances with closed forms. For example, van der Wilk et al. [41] use Gaussian processes to learn these invariances from data by maximizing the marginal likelihood. Further, the connection is often approached in the opposite direction—by looking for kernels that satisfy certain invariance properties [15, 38]. We instead approach the connection directly via data augmentation, and show that even complicated augmentation procedures akin to those used in practice can be represented as a kernel method.

3 Data Augmentation as a Kernel

To begin our study of data augmentation, we propose and investigate a model of augmentation as a Markov process, inspired by the general manner in which the process is applied—via the composition of multiple different types of transformations. Surprisingly, we show that this augmentation model combined with a k -nearest neighbor (k -NN) classifier is asymptotically equivalent to a kernel classifier, where the kernel is a function of the base transformations. While the technical details of the section can be skipped on a first reading, the central message is that kernels appear naturally in relation to data augmentation, even when we do not start with a kernel classifier. This provides additional motivation to study kernel classifiers trained on augmented data, as in Section 4.

Markov Chain Augmentation Process. In data augmentation, the aim is to perform class-preserving transformations to the original training data to improve generalization. As a concrete example, a classifier that correctly predicts an image of the number ‘1’ should be able to predict this number whether or not the image is slightly rotated, translated, or blurred. It is therefore common to pick some number of augmentations (e.g., for images: rotation, zoom, blur, flip, etc.), and to create synthetic examples by taking an original data point and applying a sequence of these augmentations. To model this process, we consider the following procedure: given a data point, we pick augmentations from a set at random, applying them one after the other. To avoid deviating too far, with some probability we discard the point and start over from a random point in the original dataset. We formalize this below.

Definition 1 (Markov chain augmentation model). Given a dataset of n examples $z_i = (x_i, y_i) \in \mathcal{X} \times \mathcal{Y}$, we *augment* the dataset via *augmentation matrices* A_1, A_2, \dots, A_m , for $A_j \in \mathbb{R}^{\Omega \times \Omega}$, which are stochastic transition matrices over a finite state space of possible labeled (augmented) examples $\Omega := \mathcal{X} \times \mathcal{Y}$. We model this via a discrete time Markov chain with the transitions:

- With probability proportional to β_j , a *transition* occurs via augmentation matrix A_j .
- With probability proportional to γ_i , a *retraction* to the training set occurs, and the state resets to z_i .

For example, the probability of retracting to training example z_1 is $\gamma_1 / (\gamma_1 + \dots + \gamma_n + \beta_1 + \dots + \beta_m)$. The augmentation process starts from any point and follows Definition 1 for an arbitrary amount of time. The retraction steps intuitively keep the final distribution grounded closer to the original training points.

From Definition 1, by conditioning on which transition is chosen, it is evident that the entire process is equivalent to a Markov chain whose transition matrix is the weighted average of the base transitions. Note that the transition matrices A_j do not need to be materialized but are implicit from the description of the augmentation. A concrete example is given in Section B.2. Without loss of generality, we let all rates β_j, γ_i be normalized with $\sum_j \gamma_i = 1$. Let $\{e_\omega\}_{\omega \in \Omega}$ be the standard basis of Ω , and let e_{z_i} be the basis element corresponding to z_i . The resulting transition matrix and stationary distribution are given below; proofs and additional details are provided in Appendix A. This describes the long-run distribution of the augmented dataset.

Proposition 1. *The described augmentation process is a Markov chain with transition matrix:*

$$R = \left(1 + \sum_{j=1}^m \beta_j\right)^{-1} \left[\sum_{i=1}^m \beta_j A_j + \sum_{i=1}^n \gamma_i (\mathbf{1} e_{z_i}^\top)\right].$$

Lemma 1 (Stationary distribution). *The stationary distribution is given by:*

$$\pi = \rho^\top (I(\beta + 1) - A)^{-1}, \text{ where } A = \sum_{j=1}^m \beta_j A_j, \beta = \sum_{j=1}^m \beta_j, \rho = \sum_{i=1}^n \gamma_i e_{z_i}. \quad (1)$$

Lemma 1 agrees intuitively with the augmentation process: When all $\beta_j \approx 0$ (i.e., low rate of augmentation), Lemma 1 implies that the stationary distribution π is close to ρ , the original data distribution. As β_j increases, the stationary distribution becomes increasingly distorted by the augmentations.

Classification Yields a Kernel. Using our proposed model of augmentation, we can show that classifying an unseen example using augmented data results in a kernel classifier. In doing so, we can observe the effect that augmentation has on the learned feature representation of the original data. We discuss several additional uses and extensions of the result itself in Appendix A.1.

Theorem 1. *Consider running the Markov chain augmentation process in Definition 1 and classifying an unseen example $x \in \mathcal{X}$ using an asymptotically Bayes-optimal classifier, such as k -nearest neighbors. Suppose that the A_i are time-reversible with equal stationary distributions. Then in the limit as time $T \rightarrow \infty$ and $k \rightarrow \infty$, this classification has the following form:*

$$\hat{y} = \text{sign} \sum_{i=1}^n y_i \alpha_{z_i} K_{x_i, x}, \quad (2)$$

where $\alpha \in \mathbb{R}^\Omega$ is supported only on the dataset z_1, \dots, z_n , and $K \in \mathbb{R}^{\Omega \times \Omega}$ is a kernel matrix (i.e., K is symmetric positive definite and non-negative) depending only on all the augmentations A_j, β_j .

Theorem 1 follows from formulating the stationary distribution (Lemma 1) as $\pi = \alpha^\top K$ for a kernel matrix K and $\alpha \in \mathbb{R}^\Omega$. Noting that k -NN asymptotically acts as a Bayes classifier, selecting the most probable label according to this stationary distribution, leads to (2).¹ In Appendix A, we include a closed form for α and K along with the proof. We include details and examples, and elaborate on the strength of the assumptions.

Takeaways. This result has two important implications: First, kernels appear naturally in relation to complex forms of augmentation, even when we do not begin with a kernel classifier. This underscores the connection between augmentation and kernels even with complicated compositional models, and also serves as motivation for our focused study on kernel classifiers in Section 4. Second, and more generally, data augmentation—a process that produces synthetic training examples from the raw data—can be understood more directly based on its effect on downstream components in the learning process, such as the features of the original data and the resulting learned model. We make this link more explicit in Section 4, and show how to exploit it in practice in Section 5.

4 Effects of Augmentation: Invariance and Regularization

In this section we build on the connection between kernels and augmentation in Section 3, exploring directly the effect of augmentation on a kernel classifier. It is commonly understood that data augmentation can be seen as a regularizer, in that it reduces generalization error but not necessarily training error [12, 42]. We make this more precise, showing that data augmentation has two specific effects: (i) increasing the invariance by averaging the features of augmented data points, and (ii) penalizing model complexity via a regularization term based on the variance of the augmented forms. These are two approaches that have been explicitly applied to get more robust performance in machine learning, though outside of the context of data augmentation. We demonstrate connections to prior work in our derivation of the feature averaging (Section 4.1) and variance regularization (Section 4.2) terms. We also validate our theory empirically (Section 4.3), and in Section 5, show the practical utility of our analysis to both kernel and deep learning pipelines.

General Augmentation Process. To illustrate the effects of augmentation, we explore it in conjunction with a general kernel classifier. In particular, suppose that we have an original kernel K with a finite-dimensional² feature map $\phi : \mathbb{R}^d \rightarrow \mathbb{R}^D$, and we aim to minimize some smooth convex loss $l : \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$

¹We use k -NN as a simple example of a nonparametric classifier, but the result holds for any asymptotically Bayes classifier.

²We focus on finite-dimensional feature maps for ease of exposition, but the analysis still holds for infinite-dimensional feature maps.

with parameter $w \in \mathbb{R}^D$ over a dataset $(x_1, y_1), \dots, (x_n, y_n)$. The original objective function to minimize is $f(w) = \frac{1}{n} \sum_{i=1}^n l(w^\top \phi(x_i); y_i)$, with two common losses being logistic $l(\hat{y}; y) = \log(1 + \exp(-y\hat{y}))$ and quadratic $l(\hat{y}; y) = (\hat{y} - y)^2$.

Now, suppose that we first augment the dataset using an augmentation kernel T . Whereas the augmentation kernel in Section 3 had a specific form based on the stationary distribution of the proposed Markov process, here we make this more general, simply requiring that for each data point x_i , $T(x_i)$ describes the distribution over data points into which x_i can be transformed. The new objective function becomes:

$$g(w) = \frac{1}{n} \sum_{i=1}^n \mathbf{E}_{t_i \sim T(x_i)} [l(w^\top \phi(t_i); y_i)] . \quad (3)$$

4.1 Data Augmentation as Feature Averaging

We begin by showing that, to first order, objective (3) can be approximated by a term that computes the average augmented feature of each data point. In particular, suppose that the applied augmentations are “local” in the sense that they do not significantly modify the feature map ϕ . Using the first-order Taylor approximation, we can expand each term around any point ϕ_0 that does not depend on t_i :

$$\mathbf{E}_{t_i \sim T(x_i)} [l(w^\top \phi(t_i); y_i)] \approx l(w^\top \phi_0; y_i) + \mathbf{E}_{t_i \sim T(x_i)} [w^\top (\phi_0 - \phi(t_i))] l'(w^\top \phi_0; y_i) .$$

Picking $\phi_0 = \mathbf{E}_{t_i \sim T(x_i)} [\phi(t_i)]$, the second term vanishes, yielding the *first-order approximation*:

$$g(w) \approx \hat{g}(w) := \frac{1}{n} \sum_{i=1}^n l(w^\top \mathbf{E}_{t_i \sim T(x_i)} [\phi(t_i)]; y_i) . \quad (4)$$

This is exactly the objective of a linear model with a new feature map $\psi(x) = \mathbf{E}_{t \sim T(x)} [\phi(t)]$, i.e., the average feature of all the transformed versions of x . If we overload notation and use $T(x, u)$ to denote the probability density of transforming x to u , this feature map corresponds to a new kernel:

$$\begin{aligned} \bar{K}(x, x') &= \langle \psi(x), \psi(x') \rangle = \langle \mathbf{E}_{u \sim T(x)} [\phi(u)], \mathbf{E}_{u' \sim T(x')} [\phi(u')] \rangle \\ &= \int_{u \in \mathbb{R}^n} \int_{u' \in \mathbb{R}^n} \langle \phi(u), \phi(u') \rangle T(x, u) T(x', u') du' du \\ &= \int_{u \in \mathbb{R}^n} \int_{u' \in \mathbb{R}^n} K(u, u') T(x, u) T(x', u') du' du \\ &= (TKT^\top)(x, x') . \end{aligned}$$

That is, training a kernel linear classifier with a particular loss function plus data augmentation is equivalent, to first order, to training a linear classifier with the same loss on an *augmented kernel* $\bar{K} = TKT^\top$, with feature map $\psi(x) = \mathbf{E}_{t \sim T(x)} [\phi(t)]$. This feature map is exactly the embedding of the distribution of transformed points around x into the reproducing kernel Hilbert space [26, 30]. This means that the first-order effect of training on augmented data is equivalent to training a support measure machine [25], with the n input distributions corresponding to the n distributions of transformed points around x_1, \dots, x_n . The new kernel \bar{K} has the effect of increasing the invariance of the model, as averaging the features from transformed inputs that are not necessarily present in the original dataset makes the features less variable to transformation.

By Jensen’s inequality, since the function l is convex, $\hat{g}(w) \leq g(w)$. In other words, if we solve the optimization problem that results from data augmentation, the resulting objective value using \bar{K} will be no larger. Further, if we assume that the loss function is strongly convex and strongly smooth, we can quantify how much the solution to the first-order approximation and the solution of the original problem with augmented data will differ (see Proposition 3 in the appendix). We validate the accuracy of this first-order approximation empirically in Section 4.3.

4.2 Data Augmentation as Variance Regularization

Next, we show that the second-order approximation of the objective on an augmented dataset is equivalent to variance regularization, making the classifier more robust. We can get an exact expression for the error by considering the second-order term in the Taylor expansion, with ζ_i denoting the remainder function from Taylor’s theorem:

$$\begin{aligned} g(w) - \hat{g}(w) &= \frac{1}{2n} \sum_{i=1}^n \mathbf{E}_{t_i \sim T(x_i)} \left[(w^\top (\phi(t_i) - \psi(x_i)))^2 l''(\zeta_i(w^\top \phi(t_i)); y_i) \right] \\ &= w^\top \left(\frac{1}{2n} \sum_{i=1}^n \mathbf{E}_{t_i \sim T(x_i)} [\Delta_{t_i, x_i} \Delta_{t_i, x_i}^\top l''(\zeta_i(w^\top \phi(t_i)); y_i)] \right) w, \end{aligned}$$

where $\Delta_{t_i, x_i} := \phi(t_i) - \psi(x_i)$ is the difference between the features of the transformed image t_i and the averaged features $\psi(x_i)$. If (as is the case for logistic and linear regression) l'' is independent of y , the error term is independent of the labels. That is, the original augmented objective g is the modified objective \hat{g} plus some regularization that is a function of the training examples, but not the labels. In other words, data augmentation has the effect of performing *data-dependent regularization*.

The second-order approximation to the objective is:

$$\tilde{g}(w) := \hat{g}(w) + \frac{1}{2n} \sum_{i=1}^n w^\top \mathbf{E}_{t_i \sim T(x_i)} [\Delta_{t_i, x_i} \Delta_{t_i, x_i}^\top] l''(w^\top \psi(x_i)) w. \quad (5)$$

For a fixed w , this error term is exactly the variance of the output $w^\top \phi(X)$, where the true data X is assumed to be sampled from the empirical data points x_i and their augmented versions specified by $T(x_i)$, weighted by $l''(w^\top \psi(x_i))$. This data-dependent regularization term favors weight vectors that produce similar outputs $w^\top \phi(x)$ and $w^\top \phi(x')$ if x' is a transformed version of x .

4.3 Validation of Approximation

We empirically validate³ the first- and second-order approximations, $\hat{g}(w)$ and $\tilde{g}(w)$, on MNIST [19] and CIFAR-10 [18] datasets, performing rotation, crop, or blur as augmentations, and using either an RBF kernel with random Fourier features [29] or LeNet (details in Appendix E.1) as a base model. Our results show that while both approximations perform reasonably well, the second-order approximation indeed results in a better approximation of the actual objective than the first-order approximation alone, validating the significance of the variance regularization component of data augmentation. In particular, in Figure 1a, we plot the difference after 10 epochs of SGD training, between the actual objective function over augmented data $g(w)$ and: (i) the first-order approximation $\hat{g}(w)$, (ii) second-order approximation $\tilde{g}(w)$, and (iii) second-order approximation without the first-order term, $f(w) + (\tilde{g}(w) - \hat{g}(w))$. As a baseline, we plot these differences relative to the difference between the augmented and non-augmented objective (i.e., the original images), $f(w)$. In Figure 1b, to see how training on approximate objectives affect the predicted test values, we plot the prediction disagreement between the model trained on true objective and the models trained on approximate objectives. Finally, Figure 1c shows that these approximations are relatively stable in terms of performance throughout the training process. For the CIFAR-10 dataset and the LeNet model (Appendix E), the results are quite similar, though we additionally observe that the first-order approximation is very close to the model trained without augmentation for LeNet, suggesting that the data-dependent regularization of the second-order term may be the dominating effect in models with learned feature maps.

4.4 Connections to Prior Work

The approximations we have provided in this section unify several seemingly disparate works.

³Code to reproduce experiments and plots is available at https://github.com/HazyResearch/augmentation_code

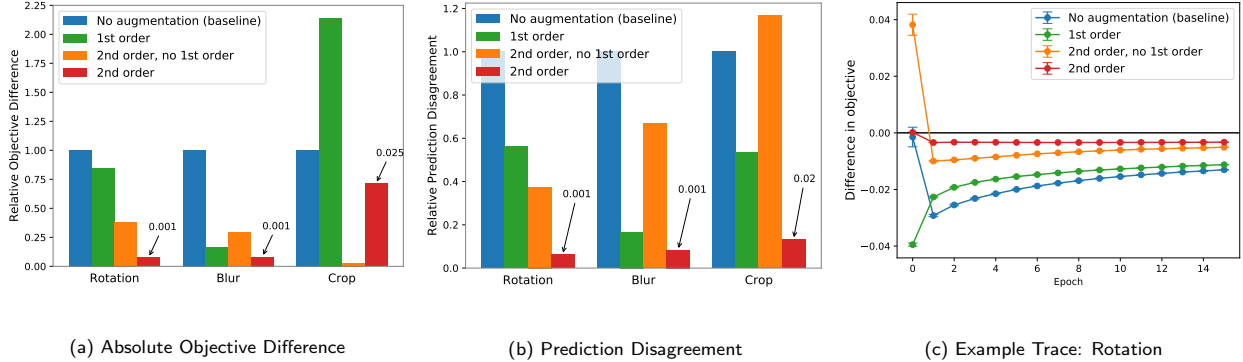


Figure 1: For the MNIST dataset, we validate that (a) the proposed approximate objectives $\hat{g}(w)$ and $\tilde{g}(w)$ are close to the true objective $g(w)$, and (b) training on the approximate objectives leads to similar predictions as training on the true objective. We plot the relative difference between the proposed approximations and the true augmented objective, in terms of difference in objective value (1a) and resulting test prediction disagreement (1b), using the non-augmented objective as a baseline. The 2nd-order approximation closely matches the true objective, particularly in terms of the resulting predictions. We observe that the accuracy of the approximations remains stable throughout training (1c). Full experiments are provided in Appendix E.

Invariant kernels. The derived first-order approximation can capture prior work in *invariant kernels* as a special case, when the transformations of interest form a group and averaging features over the group induces invariance [24, 30]. The form of the averaged kernel can then be used to learn the invariances from data [41].

Robust optimization. Our work also connects to *robust optimization*. For example, previous work [1, 3] shows that adding noise to input data has the effect of regularizing the model. Maurer & Pontil [23] bounds generalization error in terms of the empirical loss and the variance of the estimator. The second-order objective here adds a variance penalty term, thus optimizing generalization and automatically balancing bias (empirical loss) and variance with respect to the input distribution coming from the empirical data and their transformed versions (this is presumably close to the population distribution if the transforms capture the right invariance in the dataset). Though the resulting problem is generally non-convex, it can be approximated by a distributionally robust convex optimization problem, which can be efficiently solved by a stochastic procedure [28, 27].

Tangent propagation. In Section 5.3, we show that when applied to neural networks, the described second-order objective can realize classical *tangent propagation* methods [36, 37, 43] as a special case. More precisely, the second-order only term (orange in Figure 1) is equivalent to the approximation described in Zhao et al. [43], proposed there in the context of regularizing CNNs. Our results indicate that considering both the first- and second-order terms, rather than just this second-order component, in fact results in a more accurate approximation of the true objective, e.g., providing a 6–9x reduction in the resulting test prediction disagreement (Figure 1b). This suggests an approach to improve classical tangent propagation methods, explored in Section 5.3.

5 Practical Connections: Accelerating Training With Data Augmentation

We now present several proof-of-concept applications to illustrate how the theoretical insights in Section 4 can be used to accelerate training with data augmentation. First, we propose a kernel similarity metric that can be used to quickly predict the utility of potential augmentations, helping to obviate the need for guess-and-check work. Next, we explore ways to reduce training computation over augmented data, including incorporating augmentation directly in the learned features with a random Fourier features approach, and applying our derived approximation at various layers of a deep network to reduce overall computation. We

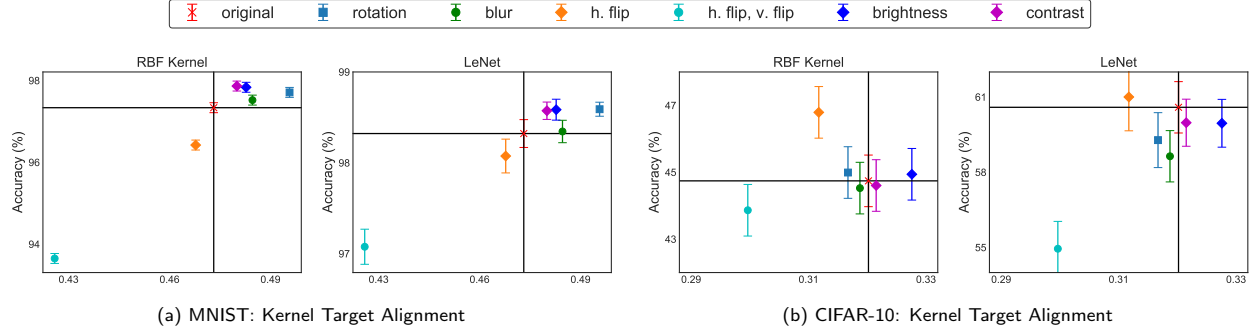


Figure 2: Accuracy vs. kernel target alignment for RBF kernel and LeNet models, for MNIST (left) and CIFAR-10 (right) datasets. This alignment metric can be used to quickly select transformations (e.g., MNIST: rotation) that improve performance and avoid bad transformations (e.g., MNIST: flips).

perform these experiments on common benchmark datasets, MNIST and CIFAR-10, as well a real-world mammography tumor-classification dataset, DDSM.

5.1 A Fast Kernel Metric for Augmentation Selection

For new tasks and datasets, manually selecting, tuning, and composing augmentations is one of the most time-consuming processes in a machine learning pipeline, yet is critical to achieving state-of-the-art performance. Here we propose a kernel alignment metric, motivated by our theoretical framework, to quickly estimate if a transformation is likely to improve generalization performance *without performing end-to-end training*.

Kernel alignment metric. Given a transformation T , and an original feature map $\phi(x)$, we can leverage our analysis in Section 4.1 to approximate the features for each data point x as $\psi(x) = \mathbf{E}_{t \sim T(x)} [\phi(t)]$. Defining the feature kernel $\bar{K}(x, x') = \psi(x)^\top \psi(x')$ and the label kernel $K_Y(y, y') = \mathbf{1}\{y = y'\}$, we can compute the *kernel target alignment* [6] between the feature kernel \bar{K} and the target kernel K_Y *without training*:

$$\hat{A}(X, \bar{K}, K_Y) = \frac{\langle \bar{K}, K_Y \rangle}{\sqrt{\langle \bar{K}, \bar{K} \rangle \langle K_Y, K_Y \rangle}},$$

where $\langle K_a, K_b \rangle = \sum_{i,j}^n K_a(x_i, x_j) K_b(x_i, x_j)$. This alignment statistic can be estimated quickly and accurately from subsamples of the data [6]. In our case, we use random Fourier features [29] as an approximate feature map $\phi(x)$ and sample $t \sim T(x)$ to estimate the averaged feature $\psi(x) = \mathbf{E}_{t \sim T(x)} [\phi(t)]$. The kernel target alignment measures the extent to which points in the same class have similar features. If this alignment is larger than that between the original feature kernel $K(x, x') = \phi(x)^\top \phi(x')$ and the target kernel, we postulate that the transformation T is likely to improve generalization. We validate this method on MNIST and CIFAR-10 with numerous transformations (rotation, blur, flip, brightness, and contrast). In Figure 2, we plot the accuracy of the kernel classifier and LeNet against the kernel target alignment. We see that there is indeed a correlation between kernel alignment and accuracy, as points tend to cluster in the upper right (higher alignment, higher accuracy) and lower left (lower alignment, lower accuracy) quadrants, indicating that this approach may be practically useful to detect the utility of a transformation prior to training.

5.2 Efficient Augmentation via Random Fourier Features

Beyond predicting the utility of an augmentation, we can also use our theory to reduce the computation required to perform augmentation on a kernel classifier—resulting, e.g., in a 4x speedup while achieving the same accuracy (MNIST, Table 1). When the transformations are affine (e.g., rotation, translation, scaling, shearing), we can perform transforms *directly on the approximate kernel features*, rather than the raw data points, thus gaining efficiency while maintaining accuracy.

Recall from Section 4 that the first-order approximation of the new feature map is given by $\psi(x) = \mathbf{E}_{t \sim T(x)} [\phi(t)]$, i.e., the average feature of all the transformed versions of x . Suppose that the transform is linear in x of the form $A_\alpha x$, where the transformation is parameterized by α . For example, a rotation by angle α has the form $T(x) = R_\alpha x$, where R_α is a $d \times d$ matrix that 2D-rotates the image x . Further, assume that the original kernel $k(x, x')$ is shift-invariant (say an RBF kernel), so that it can be approximated by random Fourier features [29]. Instead of transforming the data point x itself, we can transform the averaged feature map for x directly as:

$$\tilde{\psi}(x) = s^{-1} D^{-1/2} \left[\sum_{j=1}^s \exp(i(A_{\alpha_j}^\top \omega_1)^\top x) \quad \dots \quad \sum_{j=1}^s \exp(i(A_{\alpha_j}^\top \omega_D)^\top x) \right],$$

where $\omega_1, \dots, \omega_D$ are sampled from the spectral distribution, and $\alpha_1, \dots, \alpha_s$ are sampled from the distribution of the parameter α (e.g., uniformly from $[-15, 15]$ if the transform is rotation by α degrees). This type of random feature map has been suggested by Raj et al. [30] in the context of kernels invariant to actions of a group. Our theoretical insights in Section 4 thus connect data augmentation to invariant kernels, allowing us to leverage the approximation techniques in this area. Our framework highlights additional ways to improve this procedure: if we view augmentation as a modification of the feature map, we naturally apply this feature map to test data points as well, implicitly reducing the variance in the features of different versions of the same data point. This variance regularization is the second goal of data augmentation discussed in Section 4.

Table 1: Performance of augmented random Fourier features on MNIST, CIFAR-10, and DDSM.

Model	MNIST		CIFAR-10		DDSM	
	Acc. (%)	Time	Acc. (%)	Time	Acc. (%)	Time
No augmentation	96.1 \pm 0.1	34s	39.4 \pm 0.5	51s	57.3 \pm 6.7	27s
Traditional augmentation	97.6 \pm 0.2	220s	45.3 \pm 0.5	291s	59.4 \pm 3.2	61s
Augmented RFFs	97.6 \pm 0.1	54s	45.2 \pm 0.4	124s	58.8 \pm 5.1	34s

We validate this approach on standard image datasets MNIST and CIFAR-10, along with a real-world mammography tumor-classification dataset called Digital Database for Screening Mammography (DDSM) [17, 5, 20]. DDSM comprises 1506 labeled mammograms, to be classified as benign versus malignant tumors. In Table 1, we compare: (i) a baseline model trained on non-augmented data, (ii) a model trained on the true augmented objective, and (iii) a model that uses augmented random Fourier features. We augment via rotation between -15 and 15 degrees. All models are RBF kernel classifiers with 10,000 random Fourier features, and we report the mean accuracy and standard deviation over 10 trials. To make the problem more challenging, we also randomly rotate the test data points. The results show that augmented random Fourier features can retain 70-100% of the accuracy boost of data augmentation, with 2-4x reduction in training time.

5.3 Intermediate-Layer Feature Averaging for Deep Learning

Finally, while our theory does not hold exactly given the non-convexity of the objective, we show that our theoretical framework also suggests ways in which augmentation can be efficiently applied in deep learning pipelines. In particular, let the first k layers of a deep neural network define a feature map ϕ , and the remaining layers define a non-linear function $f(\phi(x))$. The loss on each data point is then of the form $\mathbf{E}_{t_i \sim T(x_i)} [l(f(\phi(t_i)); y_i)]$. Using the second-order Taylor expansion around $\psi(x_i) = \mathbf{E}_{t_i \sim T(x_i)} [\phi(t_i)]$, we obtain the objective:

$$\frac{1}{n} \sum_{i=1}^n l(f(\psi(x_i)); y_i) + \frac{1}{2} \mathbf{E}_{t_i \sim T(x_i)} \left[(\phi(t_i) - \psi(x_i))^\top \nabla_{\psi(x_i)}^2 l(f(\psi(x_i)); y_i) (\phi(t_i) - \psi(x_i)) \right].$$

If $f(\phi(x)) = w^\top \phi(x)$, we recover the result in Section 4 (Equation 5). Operationally, we can carry out the forward pass on all transformed versions of the data points up to layer k (i.e., computing $\phi(t_i)$), and then averaging the features and continuing with the remaining layers using this averaged feature, thus reducing computation.

We train with this approach, applying the approximation at various layers of a LeNet network using rotation as the augmentation. To get a rough measure of tradeoff between accuracy of the model and computation,

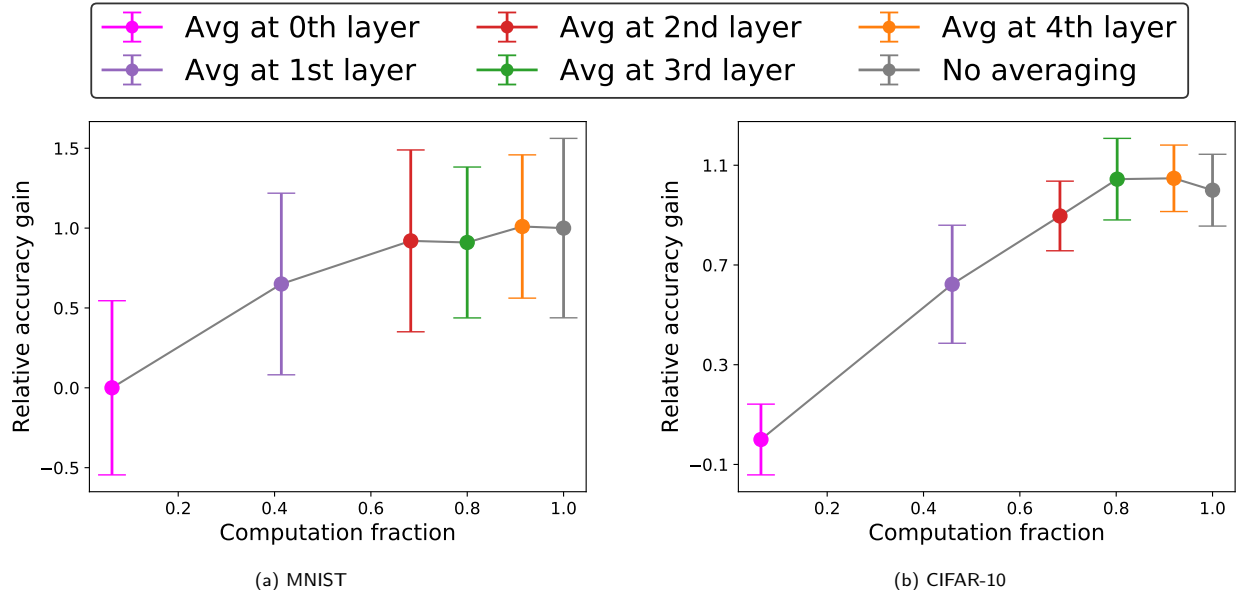


Figure 3: Accuracy gain relative to baseline (no augmentation) when averaging at various layers of a LeNet network. Approximation at earlier layers saves computation but can reduce the fidelity of the approximation.

we record the fraction of time spent at each layer in the forward pass, and use this to measure the expected reduction in computation when approximating at layer k . In Figure 3, we plot the relative accuracy gain of the classifier when trained on approximate objectives against the fraction of computation time, where 0 corresponds to accuracy (averaged over 10 trials) of training on original data, and 1 corresponds to accuracy of training on true augmented objective $g(w)$. These results indicate, e.g., that this approach can reduce computation by 30%, while maintaining 92% of the accuracy gain (red, Figure 3a). In Appendix E.4, we demonstrate similar results in terms of the test prediction distribution throughout training.

Connection to tangent propagation. If we perform the described averaging *before the very first layer* and use the analytic form of the gradient with respect to the transformations (i.e., *tangent vectors*), this procedure recovers *tangent propagation* [36]. The connection between augmentation and tangent propagation in this special case was recently observed in Zhao et al. [43]. However, as we see in Figure 3, applying the approximation at the first layer (standard tangent propagation) can in fact yield very poor accuracy results—similar to performing no augmentation—showing that our more general approximation can improve this approach in practice.

6 Conclusion

We have taken steps to establish a theoretical base for modern data augmentation. First, we analyze a general Markov process model and show that the k -nearest neighbors classifier applied to augmented data is asymptotically equivalent to a kernel classifier, illustrating the effect that augmentation has on downstream representation. Next we show that local transformations for data augmentation can be approximated by first-order feature averaging and second-order variance regularization components, having the effects of inducing invariance and reducing model complexity. We use our insights to suggest ways to accelerate training for kernel and deep learning pipelines. Generally, a tension exists between incorporating domain knowledge more naturally via data augmentation, or through more principled kernel approaches. We hope our work will enable easier translation between these two paths, leading to simpler and more theoretically grounded applications of data augmentation.

Acknowledgments

We thank Fred Sala and Madalina Fiterau for helpful discussion, and Avner May for providing detailed feedback on earlier versions.

We gratefully acknowledge the support of DARPA under Nos. FA87501720095 (D3M) and FA86501827865 (SDH), NIH under No. N000141712266 (Mobilize), NSF under Nos. CCF1763315 (Beyond Sparsity) and CCF1563078 (Volume to Velocity), ONR under No. N000141712266 (Unifying Weak Supervision), the Moore Foundation, NXP, Xilinx, LETI-CEA, Intel, Google, NEC, Toshiba, TSMC, ARM, Hitachi, BASF, Accenture, Ericsson, Qualcomm, Analog Devices, the Okawa Foundation, and American Family Insurance, and members of the Stanford DAWN project: Intel, Microsoft, Teradata, Facebook, Google, Ant Financial, NEC, SAP, and VMWare. The U.S. Government is authorized to reproduce and distribute reprints for Governmental purposes notwithstanding any copyright notation thereon. Any opinions, findings, and conclusions or recommendations expressed in this material are those of the authors and do not necessarily reflect the views, policies, or endorsements, either expressed or implied, of DARPA, NIH, ONR, or the U.S. Government.

References

- [1] Bishop, C. M. Training with noise is equivalent to tikhonov regularization. *Neural computation*, 7(1): 108–116, 1995.
- [2] Burges, C. J. C. Geometry and invariance in kernel based methods. *Advances in Kernel Methods*, pp. 89–116, 1999.
- [3] Chapelle, O., Weston, J., Bottou, L., and Vapnik, V. Vicinal risk minimization. In Leen, T. K., Dietterich, T. G., and Tresp, V. (eds.), *Advances in Neural Information Processing Systems 13*, pp. 416–422. MIT Press, 2001.
- [4] Cireřan, D. C., Meier, U., Gambardella, L. M., and Schmidhuber, J. Deep, big, simple neural nets for handwritten digit recognition. *Neural Computation*, 22(12):3207–3220, 2010.
- [5] Clark, K., Vendt, B., Smith, K., Freymann, J., Kirby, J., Koppel, P., Moore, S., Phillips, S., Maffitt, D., Pringle, M., et al. The cancer imaging archive (TCIA): maintaining and operating a public information repository. *Journal of digital imaging*, 26(6):1045–1057, 2013.
- [6] Cristianini, N., Shawe-Taylor, J., Elisseeff, A., and Kandola, J. S. On kernel-target alignment. In *Neural Information Processing Systems*, 2002.
- [7] Cubuk, E. D., Zoph, B., Mane, D., Vasudevan, V., and Le, Q. V. Autoaugment: Learning augmentation policies from data. *arXiv preprint arXiv:1805.09501*, 2018.
- [8] Decoste, D. and Schölkopf, B. Training invariant support vector machines. *Machine Learning*, 46(1): 161–190, 2002.
- [9] Demyanov, S., Bailey, J., Kotagiri, R., and Leckie, C. Invariant backpropagation: how to train a transformation-invariant neural network. *arXiv:1502.04434*, 2015.
- [10] Dosovitskiy, A., Springenberg, J. T., Riedmiller, M., and Brox, T. Discriminative unsupervised feature learning with convolutional neural networks. In *Neural Information Processing Systems*, 2014.
- [11] Dosovitskiy, A., Fischer, P., Springenberg, J. T., Riedmiller, M., and Brox, T. Discriminative unsupervised feature learning with exemplar convolutional neural networks. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 38(9):1734–1747, 2016.
- [12] Goodfellow, I., Bengio, Y., and Courville, A. *Deep Learning*. MIT Press, 2016. <http://www.deeplearningbook.org>.
- [13] Graham, B. Fractional max-pooling. *arXiv:1412.6071*, 2014.

- [14] Györfi, L., Kohler, M., Krzyzak, A., and Walk, H. *A distribution-free theory of nonparametric regression*. Springer Science & Business Media, 2006.
- [15] Haasdonk, B. and Burkhardt, H. Invariant kernel functions for pattern analysis and machine learning. *Machine Learning*, 68(1):35–61, 2007.
- [16] He, K., Zhang, X., Ren, S., and Sun, J. Identity mappings in deep residual networks. In *European Conference on Computer Vision*, 2016.
- [17] Heath, M., Bowyer, K., Kopans, D., Moore, R., and Kegelmeyer, P. The digital database for screening mammography. *Digital mammography*, pp. 431–434, 2000.
- [18] Krizhevsky, A. and Hinton, G. Learning multiple layers of features from tiny images. 2009.
- [19] LeCun, Y., Bottou, L., Bengio, Y., and Haffner, P. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- [20] Lee, R. S., Gimenez, F., Hoogi, A., and Rubin, D. Curated breast imaging subset of ddsms. *The Cancer Imaging Archive*, 2016.
- [21] Leen, T. From data distributions to regularization in invariant learning. In *Neural Information Processing Systems*, 1995.
- [22] Lu, X., Zheng, B., Velivelli, A., and Zhai, C. Enhancing text categorization with semantic-enriched representation and training data augmentation. *Journal of the American Medical Informatics Association*, 13(5):526–535, 2006.
- [23] Maurer, A. and Pontil, M. Empirical Bernstein bounds and sample variance penalization. In *Conference on Computational Learning Theory*, 2009.
- [24] Mroueh, Y., Voinea, S., and Poggio, T. A. Learning with group invariant features: A kernel perspective. In *Neural Information Processing Systems*, 2015.
- [25] Muandet, K., Fukumizu, K., Dinuzzo, F., and Schölkopf, B. Learning from distributions via support measure machines. In Pereira, F., Burges, C. J. C., Bottou, L., and Weinberger, K. Q. (eds.), *Advances in Neural Information Processing Systems 25*, pp. 10–18. Curran Associates, Inc., 2012.
- [26] Muandet, K., Fukumizu, K., Sriperumbudur, B., Schölkopf, B., et al. Kernel mean embedding of distributions: A review and beyond. *Foundations and Trends® in Machine Learning*, 10(1-2):1–141, 2017.
- [27] Namkoong, H. and Duchi, J. Stochastic gradient methods for distributionally robust optimization with f-divergences. In *Neural Information Processing Systems*, 2016.
- [28] Namkoong, H. and Duchi, J. C. Variance-based regularization with convex objectives. In *Neural Information Processing Systems*, 2017.
- [29] Rahimi, A. and Recht, B. Random features for large-scale kernel machines. In *Neural Information Processing Systems*, 2007.
- [30] Raj, A., Kumar, A., Mroueh, Y., Fletcher, P. T., and Schölkopf, B. Local group invariant representations via orbit embeddings. *International Conference on Artificial Intelligence and Statistics*, 2017.
- [31] Ratner, A. J., Ehrenberg, H., Hussain, Z., Dunnmon, J., and Ré, C. Learning to compose domain-specific transformations for data augmentation. In *Neural Information Processing Systems*, 2017.
- [32] Rifai, S., Dauphin, Y. N., Vincent, P., Bengio, Y., and Muller, X. The manifold tangent classifier. In *Neural Information Processing Systems*, 2011.
- [33] Sajjadi, M., Javanmardi, M., and Tasdizen, T. Regularization with stochastic transformations and perturbations for deep semi-supervised learning. In *Neural Information Processing Systems*, 2016.

- [34] Schölkopf, B., Burges, C., and Vapnik, V. Incorporating invariances in support vector learning machines. In *International Conference on Artificial Neural Networks*, 1996.
- [35] Sietsma, J. and Dow, R. J. Creating artificial neural networks that generalize. *Neural Networks*, 4(1): 67–79, 1991.
- [36] Simard, P., Victorri, B., LeCun, Y., and Denker, J. Tangent prop-a formalism for specifying selected invariances in an adaptive network. In *Neural Information Processing Systems*, 1992.
- [37] Simard, P., LeCun, Y., Denker, J., and Victorri, B. Transformation invariance in pattern recognition—tangent distance and tangent propagation. *Neural Networks: Tricks of the Trade*, pp. 549–550, 1998.
- [38] Teo, C. H., Globerson, A., Roweis, S. T., and Smola, A. J. Convex learning with invariances. In *Neural Information Processing Systems*, 2008.
- [39] Tibshirani, R. and Wasserman, L. Nonparametric regression and classification, 2018. URL <http://www.stat.cmu.edu/~larry/=sml/NonparametricPrediction.pdf>.
- [40] Uhlich, S., Porcu, M., Giron, F., Enenkl, M., Kemp, T., Takahashi, N., and Mitsufuji, Y. Improving music source separation based on deep neural networks through data augmentation and network blending. In *International Conference on Acoustics, Speech and Signal Processing*, 2017.
- [41] van der Wilk, M., Bauer, M., John, S., and Hensman, J. Learning invariances using the marginal likelihood. In Bengio, S., Wallach, H., Larochelle, H., Grauman, K., Cesa-Bianchi, N., and Garnett, R. (eds.), *Advances in Neural Information Processing Systems 31*, pp. 9960–9970. Curran Associates, Inc., 2018.
- [42] Zhang, C., Bengio, S., Hardt, M., Recht, B., and Vinyals, O. Understanding deep learning requires rethinking generalization. In *International Conference on Learning Representations*, 2017.
- [43] Zhao, J., Li, J., Zhao, F., Yan, S., and Feng, J. Marginalized CNN: Learning deep invariant representations. In *British Machine Vision Conference*, 2017.

A Omitted Proofs and Results From Section 3

Here we provide additional details and proofs from Section 3. First, we prove Lemma 1 characterizing the stationary distribution of the Markov chain augmentation process.

Proof of Lemma 1. Recall that the stationary distribution satisfies $\pi R = \pi$.

Under the given notation, we can express R as

$$R = \frac{A + \mathbf{1}\rho^\top}{\beta + 1}.$$

Assume for now that $I(\beta + 1) - A$ is invertible. Notice that

$$\begin{aligned} \rho^\top (I(\beta + 1) - A)^{-1} \frac{A}{\beta + 1} &= \rho^\top (I(\beta + 1) - A)^{-1} \left(\frac{A - I(\beta + 1)}{\beta + 1} + I \right) \\ &= -\frac{\rho^\top}{\beta + 1} + \rho^\top (I(\beta + 1) - A)^{-1} \end{aligned}$$

Also, $A\mathbf{1} = \sum_j \beta_j (A_j \mathbf{1}) = \beta \mathbf{1}$, so we know that $(I(\beta + 1) - A)\mathbf{1} = \mathbf{1}$. So the inverse satisfies $(I(\beta + 1) - A)^{-1} \mathbf{1} = \mathbf{1}$ as well. Thus,

$$\begin{aligned} \rho^\top (I(\beta + 1) - A)^{-1} R &= \rho^\top (I(\beta + 1) - A)^{-1} \frac{A}{\beta + 1} + \rho^\top \frac{\mathbf{1}\rho^\top}{\beta + 1} \\ &= -\frac{\rho^\top}{\beta + 1} + \rho^\top (I(\beta + 1) - A)^{-1} + \frac{\rho^\top}{\beta + 1} \\ &= \rho^\top (I(\beta + 1) - A)^{-1} \end{aligned}$$

It follows that $\pi = \rho^\top (I(\beta + 1) - A)^{-1}$ is the stationary distribution of this chain.

Finally, we show that $I(\beta + 1) - A$ is invertible as follows. By the Gershgorin Circle Theorem, the eigenvalues of A lie in the union of the discs $B(A_{ii}, \sum_{j \neq i} |A_{ij}|) = B(A_{ii}, \beta - A_{ii})$. In particular, the eigenvalues have real part bounded by β . Thus $I(\beta + 1) - A$ has all eigenvalues with real part at least 1, hence is invertible. \square

For convenience, we now restate the main Theorem of Section 3.

Theorem 1. *Consider running the Markov chain augmentation process in Definition 1 where the base augmentations preserve labels, and classifying an unseen example $x \in \mathcal{X}$ using k -nearest neighbors. Suppose that the A_i are time-reversible with equal stationary distributions. Then there are coefficients α_{z_i} and a kernel K depending only on the augmentations, such that in the limit as the number of augmented examples goes to infinity, this classification procedure is equivalent to a kernel classifier*

$$\hat{y} = \text{sign} \sum_{i=1}^n y_i \alpha_{z_i} K_{x_i, x}. \quad (6)$$

For the remainder of this section, we will refer to K alternately as a matrix $\mathbb{R}^{\Omega \times \Omega}$ or as a function $\Omega \times \Omega \rightarrow \mathbb{R}$, with corresponding notation K_{z_i, z_j} and $K(z_i, z_j)$.

Classification process. Suppose that we receive a new example $x \in \mathcal{X}$ with unknown label y . Consider running the augmentation process for time T and determining the label for x via k -nearest neighbors.⁴ Then in the limit as $T \rightarrow \infty$, we predict

$$\hat{y} = \arg \max_{y \in \mathcal{Y}} \pi((x, y)). \quad (7)$$

⁴This works for any label-preserving non-parametric model.

In other words, as the number of augmented training examples increases, k -NN approaches a Bayes classifier: it will select the class for x that is most probable in π .

We now show that under additional mild assumptions on the base augmentations A_j , applying this classification process after the Markov chain augmentation process is equivalent to a kernel classifier. In particular, *suppose that the Markov chains corresponding to the A_j are all time-reversible and there is a positive distribution π_0 that is stationary for all A_j .* This condition is not restrictive in practice, as we discuss in Section A.1. Under these assumptions, the stationary distribution can be expressed in terms of a kernel matrix.

Lemma 2. *The stationary distribution (1) can be written as $\pi = \alpha^\top K$, where the vector $\alpha \in \mathbb{R}^\Omega$ is supported only on the dataset z_1, \dots, z_n , and K is a kernel matrix (i.e., K symmetric positive definite and non-negative) depending only on the augmentations A_j, β_j .*

Proof of Lemma 2. Let $\Pi_0 = \text{diag}(\pi_0)$. The stationary distribution can be written as

$$\begin{aligned}\pi &= \rho^\top (I(\beta + 1) - A)^{-1} \\ &= \rho^\top ((\beta + 1)\Pi_0^{-1}\Pi_0 - A\Pi_0^{-1}\Pi_0)^{-1} \\ &= \rho^\top \Pi_0^{-1}((\beta + 1)\Pi_0^{-1} - A\Pi_0^{-1})^{-1}\end{aligned}$$

(Π_0 is invertible from the assumption that π_0 is supported everywhere).

Letting $\alpha = \Pi_0^{-1}\rho$ and $K = ((\beta + 1)\Pi_0^{-1} - A\Pi_0^{-1})^{-1}$, we have $\pi = \alpha^\top K$. Clearly, α is supported only on the dataset z_1, \dots, z_n since ρ is. It remains to show that K is a kernel.

The detailed balance condition of time-reversible Markov chains states that $\pi_0(u)A_j(u, v) = \pi_0(v)A_j(v, u)$ for all augmentations A_1, \dots, A_m and $u, v \in \Omega$. This can be rewritten $A_j(u, v)\pi_0(v)^{-1} = A_j(v, u)\pi_0(u)^{-1}$ or $A_j\Pi_0^{-1} = (A_j\Pi_0^{-1})^\top$, so that $A_j\Pi_0^{-1}$ is symmetric. This implies that $A\Pi_0^{-1}$ and in turn K are symmetric.

The positivity of K follows from the Gershgorin Circle Theorem, similar to the last part of the proof of Lemma 1. To show this, it suffices to show positivity of $((\beta + 1)I - A)\Pi_0^{-1} = \Pi_0^{-1}\Pi_0((\beta + 1)I - A)\Pi_0^{-1}$; we have already shown it is symmetric. Thus it suffices to show positivity of $Z = \Pi_0((\beta + 1)I - A)$.⁵ In particular, the eigenvalues of Z are in the union of the discs $D_i = B((\pi_0)_i(\beta + 1 - A_{ii}), (\pi_0)_i(\beta - A_{ii}))$.⁶ Note that D_i has real part at least $(\pi_0)_i$, and therefore the eigenvalues of Z are at least $\min_i (\pi_0)_i > 0$.

Finally, we need to show that K is a nonnegative matrix; it suffices to show this for $(I(\beta + 1) - A)^{-1}$. Note that $\frac{1}{\beta}A$ is a stochastic matrix, hence the spectral radius of $\frac{A}{\beta+1}$ is bounded by $\frac{\beta}{\beta+1} < 1$. Therefore we can expand

$$\begin{aligned}(I(\beta + 1) - A)^{-1} &= \frac{1}{\beta + 1} \left(I - \frac{A}{\beta + 1} \right)^{-1} \\ &= \frac{1}{\beta + 1} \sum_{n=0}^{\infty} \left(\frac{A}{\beta + 1} \right)^n.\end{aligned}$$

This is a sum of nonnegative matrices, so K is nonnegative. □

By Lemma 2, since α is supported only on the training set z_1, \dots, z_n , the stationary probabilities can be expressed as $\pi(z) = \sum_{u \in \Omega} \alpha(u)K(u, z) = \sum_{i=1}^n \alpha(z_i)K(z_i, z)$.

⁵This follows from the characterization of A positive definite as $x^\top Ax > 0 \forall x \neq 0$.

⁶Where $B(x, r)$ is the ball centered at x of radius r .

Expanding the classification rule (7) using this explicit representation yields

$$\begin{aligned}\hat{y} &= \arg \max_{y \in \{-1, 1\}} \sum_{i=1}^n \alpha(z_i) K((x_i, y_i), (x, y)) \\ &= \text{sign} \sum_{y \in \{-1, 1\}} y \sum_{i=1}^n \alpha(z_i) K((x_i, y_i), (x, y)).\end{aligned}$$

Finally, suppose, as is common practice, that our augmentations A_j do not change the label y . In this case, we overload the notation K so that $K(x_1, x_2) := K((x_1, 1), (x_2, 1)) = K((x_1, -1), (x_2, -1))$ ⁷, and the classification simplifies to equation (6),

the classification rule for a kernel-trick linear classifier using kernel K . Thus, *k-NN with data augmentation is asymptotically equivalent to a kernel classifier*. This completes the proof of Theorem 1.

Rate of Convergence. The rate at which the augmentation plus k -NN classifier approaches the kernel classifier can be decomposed into two parts: the rate at which the augmentation Markov chain mixes, and the rate which the k -NN classifier approaches the true function. The latter follows from standard generalization error bounds for the k -NN classifier. For example, if the kernel regression function $L(x) = \sum_{i=1}^n \alpha(z_i) K(x_i, x)$ is smooth enough (e.g. Lipschitz) on the underlying space $\mathcal{X} = \mathbb{R}^d$, then the expected difference between the two classifiers of the probability of misclassifying new examples scales as $n^{-1/(2+d)}$, where n is the number of samples (i.e. augmentation steps) [14]. Furthermore, the stationary distribution (1) can be further analyzed to yield the finite-sample distributions of the Markov chain, which is related to the power series expansion $\rho^\top (I(\beta + 1) - A)^{-1} = \rho^\top (\beta + 1)^{-1} (I + A/(\beta + 1) + \dots)$ of Equation (1). This in turn determines the mixing rate of the Markov chain, which converges to its stationary distribution exponentially fast with rate $\beta/(\beta + 1)$. More formal statements and proofs of these bounds are in Appendix A.2.

A.1 Discussion

Here we describe our modeling assumptions in more detail, as well as additional uses of our main result in Theorem 1. First, we note that the assumptions needed for Lemmas 1 and 2 hold for most transformations used in practice. For example, the condition behind Lemma 2 is satisfied by “reversible” augmentations, or any augmentation that has equal probability of sending a point $x \in \Omega$ to y as y to x : these augmentations have symmetric transition matrices, which are time-reversible and have a uniform stationary distribution. This class includes all deterministic lossless transformations, such as jittering, flips, and rotations for images. Furthermore, lossy transformations can be combined with their inverse transformation (such as ZoomOut for ZoomIn), possibly adding a probability of not transitioning, to form a symmetric augmentation.⁸

Second, our use of augmentation matrices implies a finite state space Ω , which is defined as the set of base training examples and all possible augmentations. Note that augmentations typically yield finite orbits (e.g. flip, rotation, zoom – as output pixel values are a subset of input values), which is consistent with this assumption. Furthermore, finiteness is always true in actual models due to the use of finite precision (e.g. floating point numbers).

Beyond serving as motivating connection between data augmentation, a process applied to the raw input data, and kernels, which affect the downstream feature representation, Theorem 1 also points to alternate ways to understand and optimize the augmentation pipeline. In particular, Lemma 2 provides a closed-form representation for the induced kernel in terms of the base augmentation matrices and rates, and we point out two potential ways this alternate classifier can be useful on top of the original augmentation process.

⁷The intuition is that K measures the similarity between examples in Ω in terms of how hard it is to augment one to the other, and this distance is the same whether y is 1 or -1 in the label-preserving case. Formally, a K satisfying this condition exists because the Markov chain is not irreducible, and an appropriate π_0 putting equal weights on $y = \pm 1$ can be chosen in Lemma 2. Note also that $K((x_1, 1), (x_2, -1)) = K((x_1, -1), (x_2, 1)) = 0$ by the same intuition; formally, A is a block matrix with $A((x_1, y), (x_2, -y)) = 0$ so $K = \Pi_0 (I(\beta + 1) - A)^{-1}$ has the same property

⁸For example, if a lossy transform sends $a, b, c \rightarrow c$, with transition matrix $\begin{pmatrix} 0 & 0 & 1 \\ 0 & 0 & 1 \\ 0 & 0 & 1 \end{pmatrix}$, it can be symmetrized to $\begin{pmatrix} 2/3 & 0 & 1/3 \\ 0 & 2/3 & 1/3 \\ 1/3 & 1/3 & 1/3 \end{pmatrix}$.

In Appendix B.1 we show that if the augmentations are changed, for example by tuning the rates or adding/removing a base augmentation, the kernel matrix can potentially be directly updated from the original kernel (opposed to re-sampling an augmented dataset and re-training the classifier).

Second, many parameters of the original process appear in the kernel directly. For example, in Appendix B.2 we show that in the simple case of a single additive Gaussian noise augmentation, the equivalent kernel is very close to a Gaussian kernel whose bandwidth is a function of the variance of the jitter. Additionally, in general the augmentation rates β_j all show up in the resulting kernel in a differentiable manner. Therefore instead of treating them as hyperparameters, there is potential to optimize the underlying parameters in the base augmentations, as well as the augmentation rates, through their role in a more tractable objective.

A.2 Convergence Rate

The following proposition from Györfi et al. [14] and Tibshirani & Wasserman [39] provides generalization bounds for a k -NN classifier when $k \rightarrow \infty$ and $k/n \rightarrow 0$ at a suitable rate. Treating the equivalent kernel classifier as the true function, this bounds the risk between the k -NN and kernel classifiers as a function of the number of augmented samples n .

Proposition 2. *Let \hat{C} be the k -NN classifier. Let C_0 be the asymptotically equivalent kernel classifier from Theorem 1 and assume it is L -Lipschitz.*

Letting $r(C) = \Pr_{(x,y) \sim \pi}(y \neq C(x))$ be the risk of a classifier C , then

$$r(C) - r(C_0) \leq O(L^{d/(2+d)} n^{-1/(2+d)}).$$

Next we analyze the convergence of the Markov chain by computing its distribution at time n . Define

$$\pi_n = \frac{\rho^\top}{\beta + 1} \left(\frac{A^n}{(\beta + 1)^{n-1}} + \sum_{i=0}^{n-1} \left(\frac{A}{\beta + 1} \right)^i \right).$$

We claim that π_n is the distribution of the combined Markov augmentation process at time n .

Recall that ρ^\top is a distribution over the original training data. We naturally suppose that the initial example is drawn from this distribution, so that $\pi_0 = \rho^\top$. Note that this matches the expression for π_n at $n = 0$. All that remains to show that this is the distribution of the Markov chain at time n is to prove that $\pi_{n+1} = \pi_n R$.

From the relations $\rho^\top \mathbf{1} = 1$ and $A\mathbf{1} = \beta\mathbf{1}$, we have

$$\rho^\top A^i R = \rho^\top A^i \left(\frac{A + \mathbf{1}\rho^\top}{\beta + 1} \right) = \rho^\top \frac{A^{i+1} + \beta^i \rho^\top}{\beta + 1}$$

for all $i \geq 0$.

Therefore

$$\begin{aligned}
\pi_n R &= \frac{\rho^\top}{\beta+1} \left(\frac{A^n}{(\beta+1)^{n-1}} + \sum_{i=0}^{n-1} \left(\frac{A}{\beta+1} \right)^i \right) R \\
&= \frac{\rho^\top}{\beta+1} \left(\frac{A^{n+1} + \beta^n I}{(\beta+1)^n} + \sum_{i=0}^{n-1} \frac{A^{i+1} + \beta^i I}{(\beta+1)^{i+1}} \right) \\
&= \frac{\rho^\top}{\beta+1} \left(\frac{A^{n+1}}{(\beta+1)^n} + \sum_{i=0}^{n-1} \frac{A^{i+1}}{(\beta+1)^{i+1}} + \frac{\beta^n I}{(\beta+1)^n} + \sum_{i=0}^{n-1} \frac{\beta^i I}{(\beta+1)^{i+1}} \right) \\
&= \frac{\rho^\top}{\beta+1} \left(\frac{A^{n+1}}{(\beta+1)^n} + \sum_{i=1}^{n-1} \frac{A^i}{(\beta+1)^i} + \frac{\beta^n I}{(\beta+1)^n} + \frac{I}{\beta+1} \frac{\left(1 - \left(\frac{\beta}{\beta+1}\right)^n\right)}{\left(1 - \frac{\beta}{\beta+1}\right)} \right) \\
&= \frac{\rho^\top}{\beta+1} \left(\frac{A^{n+1}}{(\beta+1)^n} + \sum_{i=1}^{n-1} \frac{A^i}{(\beta+1)^i} + I \right) \\
&= \pi_{n+1}.
\end{aligned}$$

The difference from the stationary distribution is

$$\begin{aligned}
\pi_n - \pi &= \pi_n - \frac{\rho^\top}{\beta+1} \sum_{i=0}^{\infty} \left(\frac{A}{\beta+1} \right)^i \\
&= \frac{\rho^\top}{\beta+1} \left(\frac{A^n}{(\beta+1)^{n+1}} - \sum_{i=n}^{\infty} \left(\frac{A}{\beta+1} \right)^i \right).
\end{aligned}$$

The ℓ_2 norm of this can be straightforwardly bounded, noting that $\|A\|_{op} \leq \beta$.

$$\begin{aligned}
\|\pi_n - \pi\|_2 &\leq \frac{1}{\beta+1} \left(\frac{\beta^n}{(\beta+1)^{(n+1)}} + \sum_{i=n}^{\infty} \left(\frac{\beta}{\beta+1} \right)^i \right) \\
&\leq \frac{1}{\beta+1} \left(\frac{\beta^n}{(\beta+1)^{(n+1)}} + \frac{\beta^n}{(\beta+1)^{n-1}} \right) \\
&= \left(\frac{\beta}{\beta+1} \right)^n \left(1 + \frac{1}{(\beta+1)^2} \right).
\end{aligned}$$

A bound on the total variation distance instead incurs an extra constant (in the dimension). This shows that the augmentation chain mixes exponentially fast, i.e. takes $O((\beta+1) \log(1/\varepsilon))$ samples to converge to a desired error from the stationary distribution.

B Kernel Transformations and Special Cases

B.1 Updated Kernel for Modified Augmentations

Our analysis of the kernel classifier in Lemma 2 yields a closed form in terms of the base augmentation matrices. This allows us to modify any kernel by changing the augmentations, producing a new kernel. For example, imagine that we start with a kernel K , which has corresponding augmentation operator A such that

$$K = (I(\beta+1) - A)^{-1}.$$

Suppose that we want to add an additional augmentation operator with stochastic transition matrix \hat{A} and rate $\hat{\beta}$. The resulting kernel is guaranteed to be a non-negative kernel by Lemma 2, and it can be computed

from the known K by expanding

$$\begin{aligned}
(I(\hat{\beta} + \beta + 1) - A - \hat{\beta}\hat{A})^{-1} &= (K^{-1} + \hat{\beta}I - \hat{\beta}\hat{A})^{-1} \\
&= \left(\left[I + (I - \hat{A})\hat{\beta}K \right] K^{-1} \right)^{-1} \\
&= K \left(I + (I - \hat{A})\hat{\beta}K \right)^{-1} \\
&= K \sum_{n=0}^{\infty} \hat{\beta}^n ((\hat{A} - I)K)^n.
\end{aligned}$$

B.2 Kernel Matrix for the Jitter Augmentation

In the context of Definition 1, consider performing a single augmentation A_j corresponding to adding Gaussian noise to an input vector. Although Definition 1 uses an approximated finite sample space, for this simple case we consider the original space $\mathcal{X} = \mathbb{R}^d$. The transition matrix A_1 is just the standard Gaussian kernel, $A_1(x, y) = (2\pi\sigma^2)^{-d/2} \exp(-\|x - y\|^2/(2\sigma^2))$. With rate β , the kernel matrix by Lemma 2 is

$$K = (I(1 + \beta) - \beta A)^{-1},$$

where we think of I as the identity operator on $\mathcal{X} \rightarrow \mathcal{X}$.

We define a d -dimensional Fourier Transform satisfying

$$\mathcal{F} \exp \left(-\frac{\|t\|^2}{2\sigma^2} \right) (\omega) = (\sigma^2)^{d/2} \exp \left(-\frac{\|\omega\|^2 \sigma^2}{2} \right).$$

Note that this Fourier Transform is its own inverse on Gaussian densities.

Therefore

$$\mathcal{F}K = \left[1 + \beta - \beta(2\pi)^{d/2} \exp \left(-\frac{\|\omega\|^2 \sigma^2}{2} \right) \right]^{-1}.$$

To compute the inverse transform of this, consider the function

$$\frac{1}{\alpha - \beta \exp \left(-\|t\|^2 \sigma^2/2 \right)} = \frac{1}{\alpha} \left[\frac{1}{1 - \frac{\beta}{\alpha} \exp \left(-\|t\|^2 \sigma^2/2 \right)} \right] = \frac{1}{\alpha} \left[1 + \sum_{i=1}^{\infty} \left(\frac{\beta}{\alpha} \right)^i \exp \left(\frac{-\|t\|^2 \sigma^2}{2} i \right) \right]$$

Applying the inverse Fourier Transform \mathcal{F}^{-1} , it becomes

$$\frac{1}{\alpha} \left[\delta(\omega) + \sum_{i=1}^{\infty} \left(\frac{\beta}{\alpha} \right)^i \frac{1}{(\sigma^2 i)^{d/2}} \exp \left(\frac{-\|t\|^2}{2\sigma^2 i} \right) \right].$$

Since the value of the kernel matrix only matters up to a constant, we can scale it by the first term. We also ignore the $\delta(\omega)$ term, which in the context of Theorem 1 only serves to emphasize that a test point in the training set should be classified as its known true label. Scaling by $\alpha(\alpha/\beta)\sigma^d$, we are left with

$$\begin{aligned}
&\sum_{i=1}^{\infty} \left(\frac{\beta}{\alpha} \right)^{i-1} \frac{1}{i^{d/2}} \exp \left(\frac{-\|t\|^2}{2\sigma^2 i} \right) \\
&= \exp \left(\frac{-\|t\|^2}{2\sigma^2} \right) + \sum_{i=1}^{\infty} \left(\frac{\beta}{\alpha} \right)^i \frac{1}{(i+1)^{d/2}} \exp \left(\frac{-\|t\|^2}{2\sigma^2(i+1)} \right).
\end{aligned}$$

Finally, after plugging in the corresponding values for α and β , notice that β is proportional to $(2\pi)^{-d/2}$, which causes the sum to be negligible.

C Additional Propositions for Section 4

A function f is α -strongly convex if for all x and x' , $f(x') \geq f(x) + \nabla f(x)^\top (x' - x) + (\alpha/2) \|x' - x\|^2$; the function f is β -strongly smooth if for all x and x' , $f(x') \leq f(x) + \nabla f(x)^\top (x' - x) + (\beta/2) \|x' - x\|^2$.

If we assume that the loss is strongly convex and strongly smooth, then the difference in objective functions $g(w)$ and $\hat{g}(w)$ can be bounded in terms of the squared-norm of w , and then the minimizer of the approximate objective $\hat{g}(w)$ is close to the minimizer of the true objective $g(w)$.

Proposition 3. *Assume that the loss function $l(x; y)$ is α -strongly convex and β -strongly smooth with respect to x , and that*

$$aI \preceq \frac{1}{n} \sum_{i=1}^n \mathbf{Cov}_{t_i \sim T(x_i)} (\phi(t_i)) \preceq bI, \quad \text{and}$$

$$\frac{1}{n} \sum_{i=1}^n \psi(x_i) \psi(x_i)^\top \succeq cI.$$

Letting $w^* = \arg \min g(w)$ and $\hat{w} = \arg \min \hat{g}(w)$, then

$$\frac{\alpha a}{2} \|w\|^2 \leq g(w) - \hat{g}(w) \leq \frac{\beta b}{2} \|w\|^2, \quad \text{and}$$

$$\|w^* - \hat{w}\|^2 \leq \frac{\beta b}{\alpha c} \|\hat{w}\|^2.$$

If $\alpha c \gg \beta b$ (that is, the covariance of $\phi(t_i)$ is small relative to the square of its expected value), then $\frac{\beta b}{\alpha c} \ll 1$, and so

$$\|w^* - \hat{w}\|^2 \ll \|\hat{w}\|^2.$$

This means that minimizing the first-order approximate objective \hat{g} will provide a fairly accurate parameter estimate for the objective g on the augmented dataset.

Proof of Proposition 3. By Taylor's theorem, for any random variable X over \mathbb{R} , there exists some remainder function $\zeta : \mathbb{R} \rightarrow \mathbb{R}$ such that

$$\begin{aligned} \mathbf{E} [l(X; y)] &= \mathbf{E} \left[l(\mathbf{E}[X]; y) + (X - \mathbf{E}[X]) l'(\mathbf{E}[X]; y) + \frac{1}{2} (X - \mathbf{E}[X])^2 l''(\zeta(X); y) \right] \\ &= l(\mathbf{E}[X]; y) + \frac{1}{2} \mathbf{E} [(X - \mathbf{E}[X])^2 l''(\zeta(X); y)]. \end{aligned}$$

The condition of $l(x; y)$ being α -strongly convex and β -strongly smooth means that $\alpha \leq l''(x) \leq \beta$ for any x . Thus

$$\frac{\alpha}{2} \mathbf{Var}(X) \leq \mathbf{E} [l(X)] - l(\mathbf{E}[X]) \leq \frac{\beta}{2} \mathbf{Var}(X).$$

It follows that (letting our random variable X be $w^\top \phi(t_i)$),

$$\frac{\alpha}{2} \cdot \frac{1}{n} \sum_{i=1}^n \mathbf{Var}_{t_i \sim T(x_i)} (w^\top \phi(t_i)) \leq g(w) - \hat{g}(w) \leq \frac{\beta}{2} \cdot \frac{1}{n} \sum_{i=1}^n \mathbf{Var}_{t_i \sim T(x_i)} (w^\top \phi(t_i)).$$

Because of the assumption that $aI \preceq \frac{1}{n} \sum_{i=1}^n \mathbf{Cov}_{t_i \sim T(x_i)} (\phi(t_i)) \preceq bI$,

$$\frac{\alpha a}{2} \|w\|^2 \leq g(w) - \hat{g}(w) \leq \frac{\beta b}{2} \|w\|^2.$$

We can bound the second derivative of $\hat{g}(w)$:

$$\nabla^2 \hat{g}(w) = \frac{1}{n} \sum_{i=1}^n \psi(x_i) \psi(x_i)^\top l''(w^\top \psi(x_i); y_i) \succeq \frac{\alpha}{n} \sum_{i=1}^n \psi(x_i) \psi(x_i)^\top \succeq \alpha c,$$

where we have used the assumption that $\frac{1}{n} \sum_{i=1}^n \psi(x_i) \psi(x_i)^\top \succeq cI$. Thus \hat{g} is (αc) -strongly convex.

We bound $\hat{g}(w^*) - \hat{g}(\hat{w})$:

$$\hat{g}(w^*) - \hat{g}(\hat{w}) = \hat{g}(w^*) - g(w^*) + g(w^*) - g(\hat{w}) + g(\hat{w}) - \hat{g}(\hat{w}) \leq 0 + 0 + \frac{\beta b}{2} \|\hat{w}\|^2,$$

where we have used the fact that $\hat{g}(w) \leq g(w)$ for all w and that w^* minimizes $g(w)$. But \hat{g} is (αc) -strongly convex, so $\hat{g}(w^*) - \hat{g}(\hat{w}) \geq \alpha c/2 \|w^* - \hat{w}\|^2$. Combining these inequalities yields

$$\|w^* - \hat{w}\|^2 \leq \frac{\beta b}{\alpha c} \|\hat{w}\|^2.$$

□

D Variance Regularization Terms for Common Loss Functions

Here we derive the variance regularization term for common loss functions such as logistic regression and multinomial logistic regression.

For logistic regression,

$$\begin{aligned} l(x; y) &= \log(1 + \exp(-yx)) \\ &= -\frac{yx}{2} + \log\left(\exp\left(\frac{yx}{2}\right) + \exp\left(-\frac{yx}{2}\right)\right) \\ &= -\frac{yx}{2} + \log\left(2 \cosh\left(\frac{yx}{2}\right)\right) \\ &= -\frac{yx}{2} + \log 2 + \log \cosh\left(\frac{yx}{2}\right). \end{aligned}$$

And so

$$\begin{aligned} l'(x; y) &= -\frac{y}{2} + \frac{\sinh\left(\frac{yx}{2}\right)}{\cosh\left(\frac{yx}{2}\right)} \cdot \frac{y}{2} \\ &= -\frac{y}{2} + \frac{y}{2} \tanh\left(\frac{yx}{2}\right) \end{aligned}$$

and

$$\begin{aligned} l''(x; y) &= \frac{y^2}{4} \operatorname{sech}^2\left(\frac{yx}{2}\right) \\ &= \frac{1}{4} \operatorname{sech}^2\left(\frac{x}{2}\right), \end{aligned}$$

since $y \in \{-1, 1\}$ and so $y^2 = 1$. Therefore,

$$g(w) - \hat{g}(w) = w^\top \left(\frac{1}{4n} \sum_{i=1}^n \mathbf{E}_{t_i \sim T(x_i)} \left[(\phi(t_i) - \psi(x_i)) (\phi(t_i) - \psi(x_i))^\top \operatorname{sech}^2\left(\frac{1}{2} \zeta_i(w^\top \phi(t_i))\right) \right] \right) w.$$

To second order, this is

$$\begin{aligned} g(w) - \hat{g}(w) &\approx w^\top \left(\frac{1}{4n} \sum_{i=1}^n \mathbf{E}_{t_i \sim T(x_i)} \left[(\phi(t_i) - \psi(x_i)) (\phi(t_i) - \psi(x_i))^\top \operatorname{sech}^2\left(\frac{1}{2} (w^\top \psi(x_i))\right) \right] \right) w \\ &= w^\top \left(\frac{1}{4n} \sum_{i=1}^n \mathbf{Cov}_{t_i \sim T(x_i)} \left(\phi(t_i) \operatorname{sech}\left(\frac{1}{2} w^\top \psi(x_i)\right) \right) \right) w. \end{aligned}$$

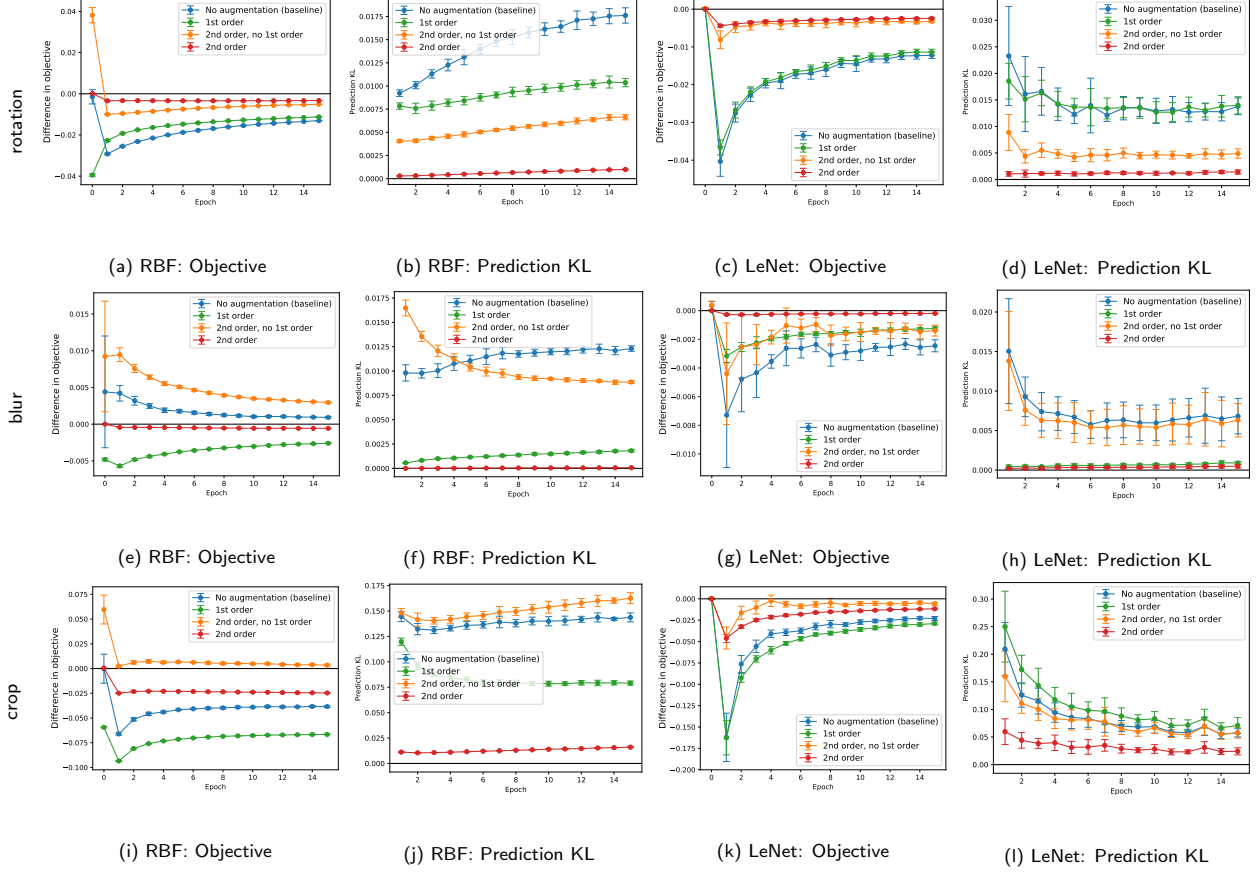


Figure 4: **MNIST Dataset.** The difference in objective value (a,e,i,c,g,k) and prediction distribution (b,f,j,d,h,l) (as measured via the KL divergence) between approximate and true objectives. In all plots, the second-order approximation tends to closely match the true objective, and to be closer than the first-order approximation or second-order component alone.

For multinomial logistic regression, we use the cross entropy loss. With the softmax probability $p_i = \frac{\exp(x_i)}{\sum \exp(x_j)}$,

$$l(x; y) = -(x_y - \log \sum \exp(x_j)).$$

The first derivative is:

$$\nabla l(x; y) = \frac{\exp(x_i)}{\sum \exp(x_j)} - \mathbf{1}_{i=y} = p - \mathbf{1}_{i=y}.$$

The second derivative is:

$$\nabla^2 l(x; y) = \text{diag}(p) - pp^T,$$

which does not depend on y .

E Additional Experiment Details and Results

E.1 First- and Second-order Approximations

For all experiments, we use the MNIST and CIFAR-10 dataset and test three representative augmentations: rotations between -15 and 15 degrees, random crops of up to 64% of the image area, and Gaussian blur. We explore our approximation for kernel classifier models, using either an RBF kernel with 10000 random Fourier

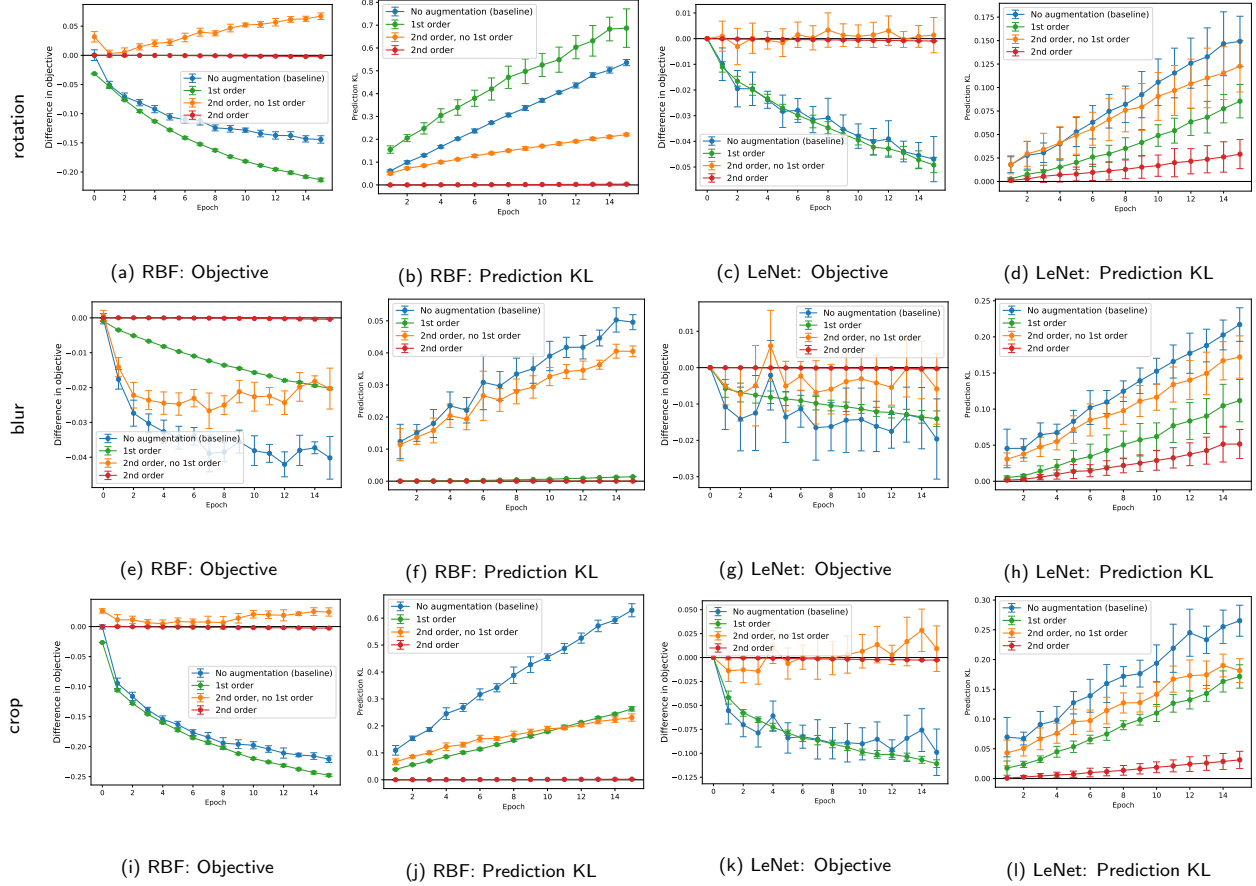


Figure 5: **CIFAR-10 Dataset.** The difference in objective value (a,e,i,c,g,k) and prediction distribution (b,f,j,d,h,l) (as measured via the KL divergence) between approximate and true objectives. In all plots, the second-order approximation tends to closely match the true objective, and to be closer than the first-order approximation or second-order component alone.

features [29] or a learned LeNet neural network [19] as our base feature map. The RBF kernel bandwidth is chosen by computing the kernel alignment, as in Section 5.1. We optimize the models using stochastic gradient descent, with learning rate 0.01, momentum 0.9, batch size 256, running for 15 epochs. We explicitly transform the images and add them to the dataset.

We validate two claims: (i) the approximate objectives are close to the true objective, and (ii) training on approximate objectives gives similar models to training on the true objective.

We plot the mean and standard deviation over 10 runs in Figure 4 and Figure 5. For claim (i), we plot the objective difference, throughout the process of training a model on the true objective, in Figure 4 (a, c, e, g, i, k) and Figure 5 (a, c, e, g, i, k). Objective difference closer to 0 is better. The second-order approximation is better than the first-order and the second-order without first-order approximation. For claim (ii), we train 5 models, each on different objectives: true objective, first-order approximation, second-order approximation, second-order without first-order approximation, and no augmentation objective. We measure the KL divergence between the predictions given by the approximate models and the predictions given by the model with true objective, as they are trained. Lower KL divergence means the prediction distributions of the approximate model is more similar to the predictions made by the true model. Figure 4 (b, d, f, h, j, l) and Figure 5 (b, d, f, h, j, l) show that the approximate models trained on first-order approximation and second-order approximation yield similar predictions to the model trained on the true objective, with the second-order approximate model being particularly close to the true model.

E.2 Kernel Alignment

For the experiment in Section 5.1, we use the same RBF kernel with 10000 random Fourier features and LeNet, as in Section E.1. We compute the kernel target alignment by collecting statistics from mini-batches of the dataset, iterating over the dataset 50 times. For the MNIST dataset, we consider rotation (between -15 and 15 degrees), Gaussian blur, horizontal flip, horizontal & vertical flip, brightness adjustment (from 0.75 to 1.25 brightness factor), contrast adjustment (from 0.65 to 1.35 contrast factor). For the CIFAR-10 dataset, we consider rotation (between -5 and 5 degrees), Gaussian blur, horizontal flip, horizontal & vertical flip, brightness adjustment (from 0.75 to 1.25 brightness factor), contrast adjustment (from 0.65 to 1.35 contrast factor). Accuracy on the validation set is obtained from SGD training over 15 epochs, with the same hyperparameters as in Section E.1, averaged over 10 trials.

E.3 Augmented Random Fourier Features

We use two standard image classification datasets MNIST and CIFAR-10, and a real-world mammography tumor-classification dataset called Digital Database for Screening Mammography (DDSM) [17, 5, 20]. DDSM comprises 1506 labeled mammograms, to be classified as benign versus malignant tumors. The DDSM images are gray-scale and of size 224×224 , which we resize to 64×64 .

We use the same RBF kernels with 10000 random Fourier features as in Section E.1. We apply rotations between -15 and 15 degrees for MNIST and CIFAR-10, and rotations between 0 and 360 degrees for DDSM, since the tumor images are rotationally invariant. We sample $s = 16$ angles to construct the augmented random Fourier feature map for MNIST and CIFAR, and $s = 36$ angles for DDSM.

To make the MNIST and CIFAR classification tasks more challenging, we also rotate the images in the test set between -15 and 15 degrees for MNIST, and between -5 and 5 degrees for CIFAR-10.

E.4 Feature Averaging for Deep Learning

In addition to the accuracy results in Section 5.3, we also explore the effect of feature averaging on prediction disagreement throughout training. We plot the difference in generalization between approximation and true objectives for LeNet in Figure 6, for rotation between -15 and 15 degrees. We again observe that approximation at earlier layers saves computation but can reduce the fidelity of the approximation.

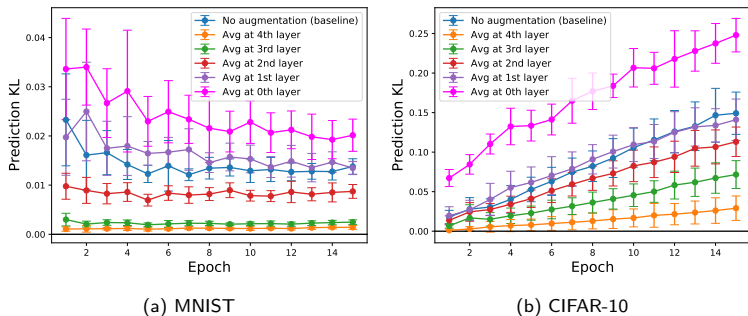


Figure 6: Difference in generalization between approximate and true objectives for LeNet in terms of KL divergence in test predictions, for MNIST (a) and CIFAR-10 (b) datasets.

E.5 Layerwise Feature Invariance in a ResNet

Finally, to provide additional motivation for the deep learning experiments in Section 5, where our theory breaks down due to non-convexity, we explore invariance with respect to deep neural networks. For each

layer l of a deep neural network, we examine the average difference in feature values when data points x are transformed according to a certain augmentation distribution T , using a model trained with data augmentation T' which has feature layers $\phi_l^{T'}$:

$$\Delta_{l,T,T'} = \sum_{i=1}^n \mathbf{E}_{z \sim T(x_i)} \left[\frac{1}{|\phi_l^{T'}(x_i)|} \|\phi_l^{T'}(x_i) - \phi_l^{T'}(z)\|^2 \right] \quad (8)$$

Specifically, in Figure 7, we examine the ratio of this measure of invariance for a model trained with data augmentation using T , and trained without any data augmentation, $\Delta_{l,T,T}/\Delta_{l,T,\emptyset}$, to see if and how training with a specific augmentation makes the layers of the network more invariant to it. We use a standard ResNet as in He et al. [16] with three blocks of nine residual units (separated by vertical dashed lines), an initial convolution layer, and a final global average pooling layer, implemented in TensorFlow⁹, trained on CIFAR-10 and averaged over ten trials. We see that training with an augmentation indeed makes the feature layers of the network more invariant to this augmentation, with the steepest change in the earlier layers (first residual block), and again in the final layer when features are pooled and averaged.

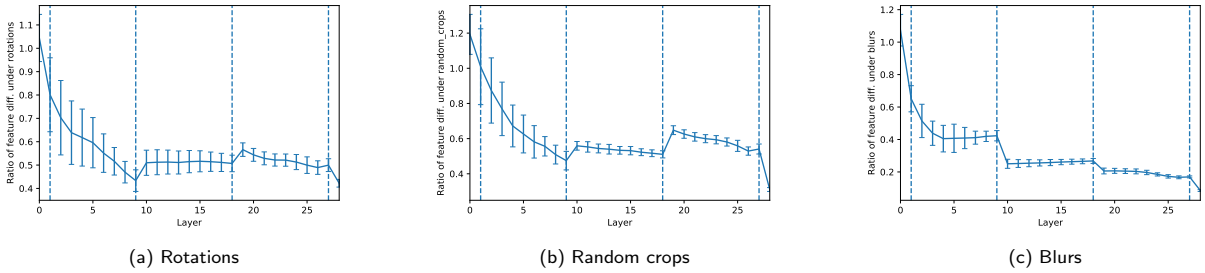


Figure 7: The ratio of average difference in features under an augmentation (8) when the network is trained with that augmentation, to when it is trained without any data augmentation, using a ResNet trained on CIFAR-10, averaged over ten trials. We see that in all but the first one or two layers, the network trained with the augmentation is indeed more invariant to it, with the steepest increase in invariance in the first block of layers and in the last global average pooling layer.

⁹<https://github.com/tensorflow/models/tree/master/official/resnet>